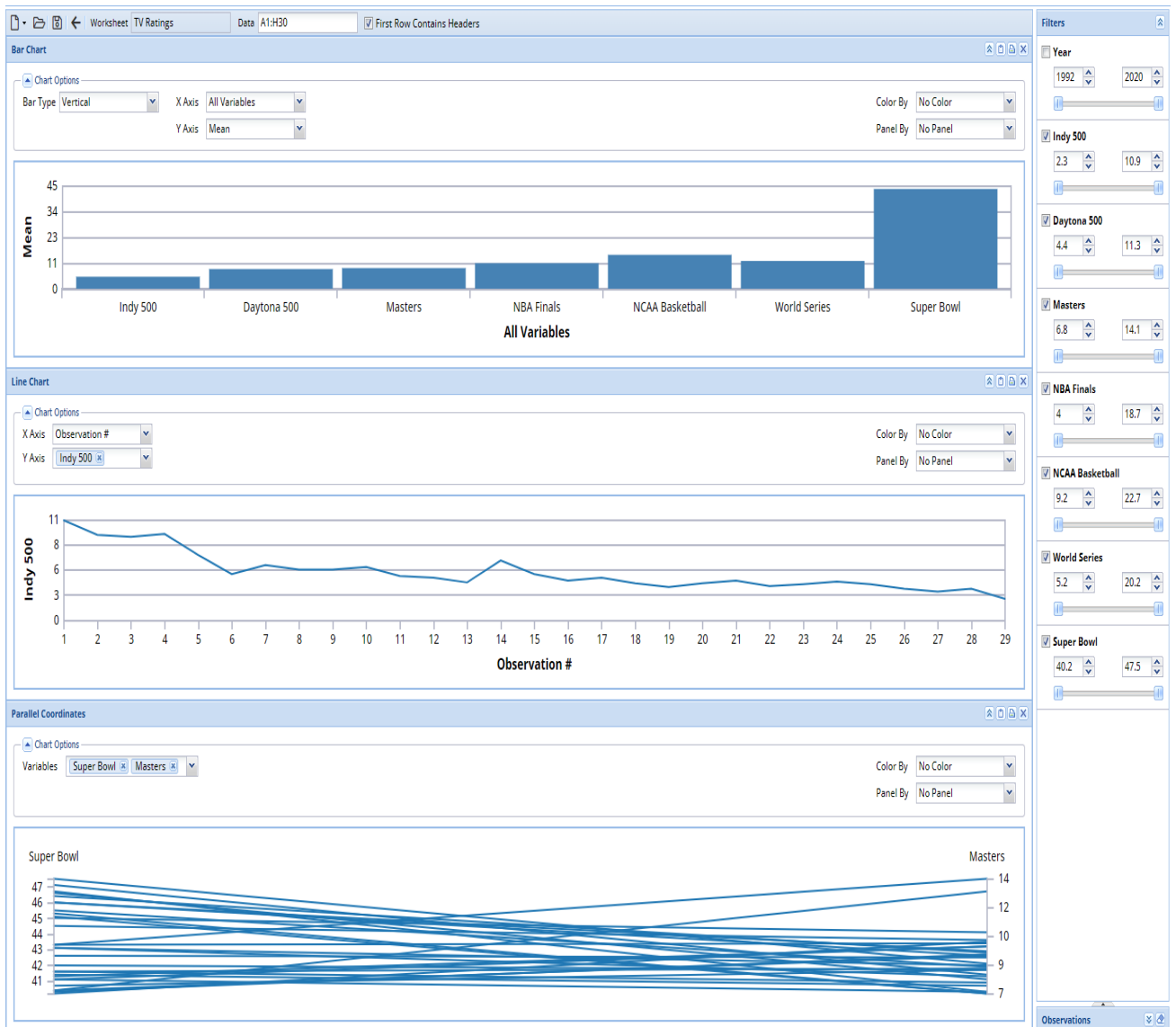


Analytic Solver

Data Science User Guide



Copyright

Software copyright 1991-2025 by Frontline Systems, Inc.

User Guide copyright 2025 by Frontline Systems, Inc.

GRG/LSGRG Solver: Portions copyright 1989 by Optimal Methods, Inc. SOCP Barrier Solver: Portions copyright 2002 by Masakazu Muramatsu. LP/QP Solver: Portions copyright 2000-2010 by International Business Machines Corp. and others. Neither the Software nor this User Guide may be copied, photocopied, reproduced, translated, or reduced to any electronic medium or machine-readable form without the express written consent of Frontline Systems, Inc., except as permitted by the Software License agreement below.

Trademarks

Frontline Solvers®, XLMiner®, Analytic Solver®, Risk Solver®, Premium Solver®, Solver SDK®, and RASON® are trademarks of Frontline Systems, Inc. Windows and Excel are trademarks of Microsoft Corp. Gurobi is a trademark of Gurobi Optimization, Inc. Knitro is a trademark of Artelys. MOSEK is a trademark of MOSEK ApS. OptQuest is a trademark of OptTek Systems, Inc. Xpress^{MP} is a trademark of FICO, Inc.

Patent Pending

Systems and Methods for Automated Risk Analysis of Machine Learning Models.

Acknowledgements

Thanks to Dan Fylstra and the Frontline Systems development team for a 25-year cumulative effort to build the best possible optimization and simulation software for Microsoft Excel. Thanks to Frontline's customers who have built many thousands of successful applications, and have given us many suggestions for improvements.

Risk Solver Pro and Risk Solver Platform have benefited from reviews, critiques, and suggestions from several risk analysis experts:

- Sam Savage (Stanford Univ. and AnalyCorp Inc.) for Probability Management concepts including SIPs, SLURPs, DISTs, and Certified Distributions.
- Sam Sugiyama (EC Risk USA & Europe LLC) for evaluation of advanced distributions, correlations, and alternate parameters for continuous distributions.
- Savvakis C. Savvides for global bounds, censor bounds, base case values, the Normal Skewed distribution and new risk measures.

How to Order

Contact Frontline Systems, Inc., P.O. Box 4288, Incline Village, NV 89450.

Tel (775) 831-0300 Fax (775) 831-0314 Email info@solver.com Web <http://www.solver.com>

Table of Contents

Table of Contents	3
Start Here: Data Science Essentials in V2025 Q1	8
Getting the Most from This User Guide	8
Desktop and Cloud Versions	8
Installing the Software	8
Understanding License and Upgrade Options	8
Getting Help Quickly	8
Finding the Examples	9
Using Existing Models	9
Getting and Interpreting Results	9
Installation and Add-Ins	10
What You Need	10
Installing the Software	10
Installing Analytic Solver Cloud	10
Installing Analytic Solver Desktop	12
Logging in the First Time	17
Uninstalling the Software	18
Activating and Deactivating the Software	18
Excel 2019, 2016 and 2013	18
Analytic Solver Data Science Overview	20
Overview	20
V2016 Release	21
V2016-R2 Release	21
V2016-R3 Release	21
V2017 Release	22
V2017-R2 Release	22
V2018 Release	23
V2019 Release	24
V2020 Release	24
What's New in Analytic Solver V2020.5	25
Easily Deploy Your Model as a Cloud Service	25
More About RASON Decision Services	26
How You Can Use RASON	27
New Time Series Simulation Functions	28
New Optimization Result Functions	28
What's New in Analytic Solver V2021	28
Lambda, Let and Box Functions	29
What's New in Analytic Solver V2021.5	29
Automate Data Mining with Find Best Model	29
Better Simulation Models with Metalog Distributions and Fitting	29
Share Data Mining and Probability Models via RASON	30

And of Course, Optimization Enhancements.....	30
Analytic Solver Upgrade.....	30
Analytic Solver Academy – Replaces Analytic Solver Basic.....	31
What's New in Analytic Solver V2022.....	31
Faster Interaction, Faster Solves.....	31
What's New in Analytic Solver V2023.....	31
Automated Risk Analysis of Machine Learning Models.....	31
Synthetic Data Generation ‘For Free’.....	32
What’s New in Analytic Solver V2023 Q1.....	32
Faster LP/Quadratic Solver and Large-Scale LP/QP Solver Engine.....	32
More Plug-in Solver Engine Improvements.....	33
Greatly Improved “Deploy Your Model to Teams” Capability.....	33
Risk Analysis of Machine Learning Models Created in Other Software.....	33
What’s New in Analytic Solver V2023 Q3.....	34
AI Agent: Ask for Help from ChatGPT “Trained on Analytic Solver”.....	34
Identify Inputs: Easily Set Up your Model for Data Updates and New Solves .	35
What’s New in Analytic Solver V2024 Q2.....	35
Solving Outside Excel.....	35
New Versions of Gurobi and OptQuest Solvers.....	36
PSI Interpreter Improvements.....	36
More About the PSI Interpreter: The Jacobian.....	37
Beyond the Jacobian: The Hessian.....	37
And There's (Much) More.....	38
What’s New in Analytic Solver V2024 Q3.....	38
New AI Assistant.....	38
Optimization Improvements.....	38
Simulation Improvements.....	39
What’s New in Analytic Solver V2025 Q1.....	39
Usability: New Flexible Dialogs.....	39
Performance: New Nonlinear Solver Power.....	39
Standard Excel Formulas, New Solving Speed.....	40
Analytic Solver Product Line.....	41
Desktop and Cloud versions.....	41
Analytic Solver Academy.....	42
Analytic Solver Upgrade.....	42
Analytic Solver Optimization.....	43
Analytic Solver Simulation.....	43
Analytic Solver Data Science.....	43
Analytic Solver Comprehensive.....	44
Data Science Ribbon Overview.....	46
Model.....	48
Get Data.....	48
Data Analysis.....	49
Time Series Analysis.....	50
Data Science.....	51
Scoring.....	53

Using Help, Licensing and Product Subsets 54

Introduction.....	54
Working with Licenses in V2025 Q1.....	54
Frontline License Manager.....	54
Managing Your Licenses.....	55
Product Selection Wizard.....	57
Getting Help.....	59
Ask Question.....	59

AI Assist.....	61
Maximizing Benefits of AI Agent and AI Assist.....	68
Help Center	69
Accessing Resources.....	69
User Guides.....	70
Example Models.....	70
Knowledge Base.....	70
Operating Mode.....	70
Support Mode.....	70
Submit a Support Ticket.....	71
Solver Academy	71
Video Tutorials/Live Webinars	72
Learn more!.....	72
Help Menu.....	72
Creating & Deploying a Workflow	73
Introduction.....	73
Creating a Workflow	73
Workflow Tab Options	76
Deploying Your Workflow	77
Posting Workflow to RASON Cloud Services	78
Deploying Your Model	80
Manually Creating a Workflow.....	81
Making/Breaking a Connection	85
Running a Workflow with a New Dataset.....	85
Multiple Workflows.....	86
Changing Options Settings.....	87
Workflow Groups	87
Deploying a Fitted Model to RASON, Power BI or Tableau	89
Introduction.....	89
RASON Deployment Menu	90
Deployment Wizard.....	91
RASON Deployment Menu	91
Conversion Exceptions	93
Manage Models	93
Cloud Service.....	94
Cloud Service Fitted Model.....	94
Power BI.....	96
Power BI Fitted Model.....	96
Tableau	100
Tableau Fitted Model.....	101
Bringing Big Data into Excel Using Apache Spark	109
Introduction.....	109
Sampling and Summarizing Big Data.....	109
Connecting to an Apache Spark Cluster.....	110
Storage Sources and Data Formats	110
Sampling from a Large Dataset	110
Summarizing a Large Dataset.....	114
Fitting a model using Feature Selection	119
What is Feature Selection?.....	119

Feature Selection Example.....	120
Text Mining	133
Introduction.....	133
Text Mining Example.....	134
Importing from a File Folder.....	134
Using Text Miner.....	137
Output Results.....	146
Classification with Concept Document Matrix.....	150
Exploring a Time Series Dataset	154
Introduction.....	154
Autocorrelation (ACF).....	154
Partial Autocorrelation Function (PACF).....	155
ARIMA.....	155
Partitioning.....	156
Examples for Time Series Analysis.....	156
Automated Risk Analysis of Machine Learning Models	164
Introduction.....	164
Risk Analysis of ML Model in Predicting Loan Defaults.....	164
Preliminary Illustration to Show Synthetic Data Compared to Known Data...	165
Illustration: Automated Risk Analysis of Just-Trained Machine Learning Model	169
.....	169
Fitting the Best Model	177
Introduction.....	177
Find Best Model Classification Example.....	177
Opening the Dataset.....	178
Partitioning the Dataset.....	178
Running Find Best Model.....	179
Interpreting the Results.....	186
Scoring New Data.....	201
Find Best Model Regression Model.....	202
Opening the Dataset.....	203
Partitioning the Dataset.....	204
Running Find Best Model.....	204
Simulation Tab.....	210
Interpreting the Results.....	210
Scoring New Data.....	226
Classifying the Iris Dataset	228
Introduction.....	228
Creating the Classification Model.....	228
Predicting Housing Prices using Multiple Linear Regression	240
Introduction.....	240
Multiple Linear Regression Example.....	240
Input.....	240
Output.....	246

Scoring New Data	262
Introduction.....	262
Scoring New Data	262
Scoring New Data Example	262
Using Data Science Psi Functions in Excel.....	272
Scoring Data Using Psi Functions	273
PsiPredict().....	274
PsiPosteriors()	276
PsiTransform().....	277
Time Series Forecasting.....	278
Time Series Forecasting.....	278
Time Series Forecasting.....	279
PsiForecast().....	283
Time Series Simulation.....	284
Scoring to a Database	289
Scoring to a Worksheet.....	297

Start Here: Data Science Essentials in V2025 Q1

Getting the Most from This User Guide

Desktop and Cloud Versions

Analytic Solver V2025 Q1 comes in two versions: **Analytic Solver Desktop** – a traditional “COM add-in” that works only in Microsoft Excel for Windows PCs (desktops and laptops), and **Analytic Solver Cloud** – a modern “Office add-in” that works in Excel for Windows and Excel for Macintosh (desktops and laptops), and also in Excel for the Web using Web browsers such as Chrome, FireFox and Safari. Your license gives you access to both versions, and your Excel workbooks and optimization, simulation and data science models work in both versions, no matter where you save them (though OneDrive is most convenient).

Installing the Software

Read the chapter “Installation and Add-Ins” for complete information on installing Analytic Solver Cloud and (if you wish) Analytic Solver Desktop. This chapter also explains how the Cloud and Desktop versions interact when both are installed, and how to install and uninstall both versions.

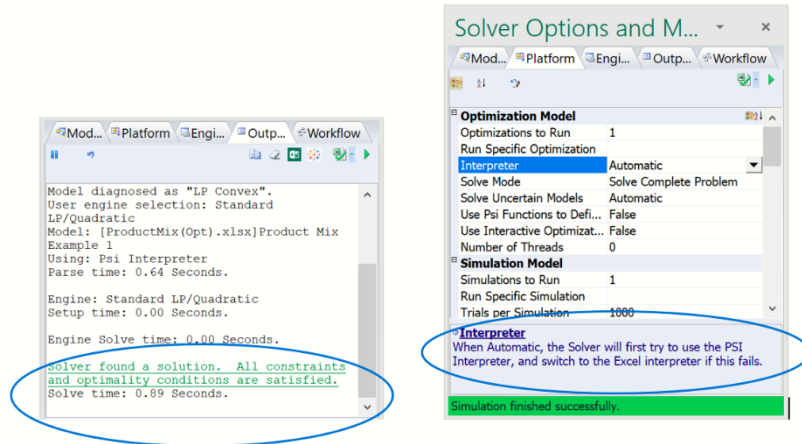
In brief, to add **Analytic Solver Cloud** version to your copy of Excel, you use the Excel Ribbon option **Insert – Get Add-ins** – no Setup program download or installation is required. To install **Analytic Solver Desktop** on Windows PCs, visit www.solver.com, login using the same email and password you used to register on Solver.com, download and run the SolverSetup program.

Understanding License and Upgrade Options

Frontline Solvers V2025 Q1 features product line that gives you access to **all** features, **all** the time for models of all sizes, and a new licensing system, tied to **you** and usable on more than one computer. Read about this in the chapter “Help, Support, Licenses and Product Versions.”

Getting Help Quickly

Choose Help on the Ribbon. You’ll see several options, starting with **Help – Help Center**. Support Live Chat, Example Models, and User Guides are also available here. In Analytic Solver Desktop (only) you can also get quick online Help by clicking any underlined caption or message in the Task Pane.



Finding the Examples

Use **Help – Examples** on the Analytic Solver or Data Science Ribbon to open a list of example optimization and simulation models, and example data sets for data science, that you can open by clicking hyperlinks. See the chapter “Help, Support, Licenses and Product Subsets” for details. Some of these examples are used and described in the **Examples** chapters.

Using Existing Models

Models created using XLMiner 4.0 and earlier can be used in Analytic Solver Data Science without any required changes.

Getting and Interpreting Results

Learn how to interpret Analytic Solver Data Science’s result messages, error messages, reports and charts using the Help file imbedded within the software. Simply go to **Help – User Guides** on the Data Science ribbon.

Installation and Add-Ins

What You Need

You can use Analytic Solver Cloud in Excel for the Web (formerly Excel Online) through a Web browser (such as Edge, Chrome, Firefox or Safari), without installing anything else. This is the simplest and most flexible option, but it requires a constant Internet connection.

To use Analytic Solver Cloud in Excel Desktop on a PC or Mac, you must have a current version of Windows or iOS installed, and **you will need the latest Excel version installed via your Office 365 subscription** – older non-subscription versions, even Excel 2019, do not have all the features and APIs needed for modern JavaScript add-ins like Analytic Solver Cloud.

To use Analytic Solver Desktop (Windows PCs only), you must have first installed Microsoft Excel 2013, 2016, 2019, or the latest Office 365 version on Windows 10, Windows 8, Windows 7, or Windows Server 2019, 2016 or 2012. (Windows Vista or Windows Server 2008 may work but are no longer supported.). It's not essential to have the standard Excel Solver installed.

Installing the Software

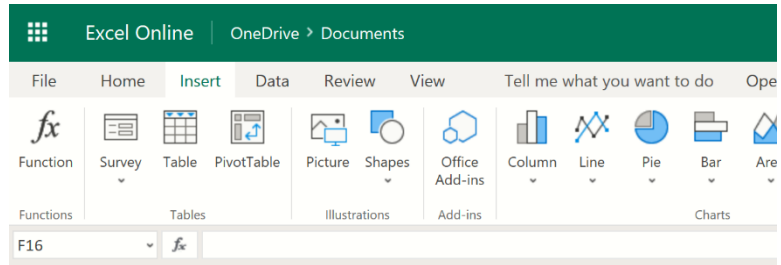
Installing Analytic Solver Cloud

Analytic Solver V2025 Q1 includes our next-generation offering, Analytic Solver Cloud – usable in the latest versions of desktop Excel for Windows and Macintosh, and in Excel for the Web. Analytic Solver Cloud is divided into **two** add-ins that work closely together (since a JavaScript add-in currently can have only one Ribbon tab): the **Analytic Solver** add-in builds optimization, simulation and decision table models, and the **Data Science** add-in builds data science or forecasting models.

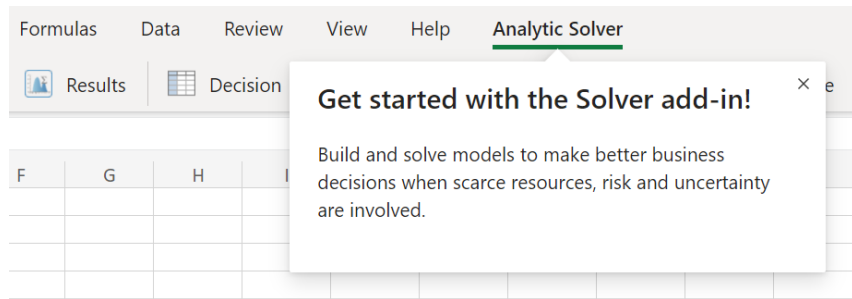
Both the Analytic Solver and Data Science add-ins support existing models created in previous versions of Analytic Solver. Your license for Analytic Solver will allow you to use Analytic Solver Desktop in desktop Excel or Analytic Solver Cloud in either desktop Excel (latest version) *or* Excel for the Web.

Once you do this, the Analytic Solver and Data Science tabs will appear on the Ribbon in each new workbook you use. To use the Analytic Solver and Data Science add-ins, **you must first “insert” them** for use in your licensed copy of desktop Excel or Excel for the Web, while you are logged into your Office 365 account.

To insert the add-ins for the first time, open desktop Excel (latest version) or Excel for the Web, click the **Insert** tab on the Ribbon, then click the button **Office Add-ins** or (if you see it) the smaller button **Get Add-ins**.

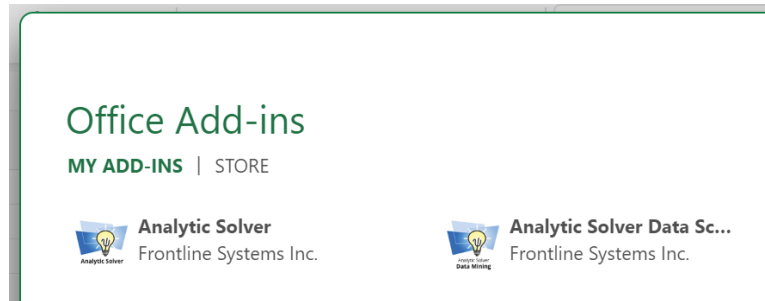


In the dialog box that appears, click the **Store** tab and type “Analytic Solver” into the Search box. Once you find the Analytic Solver add-in, click **Add**. After a moment, you should see the **Analytic Solver** tab appear on the Ribbon, with a note about how to “Get started with the Solver add-in!”, as shown below.



Repeat these steps to search for, locate and **Add** the Analytic Solver Data Science add-in. After a moment, you should see the **Data Science** tab appear on the Ribbon, with a similar “Get started” note.

After you perform these steps (one time) to insert the Analytic Solver and Analytic Solver Data Science add-ins, they will appear under "My Add-ins". If you ever need to remove the add-ins, click the “...” symbol to the right of the add-in name, then click the **Remove** choice on the dropdown menu that appears.



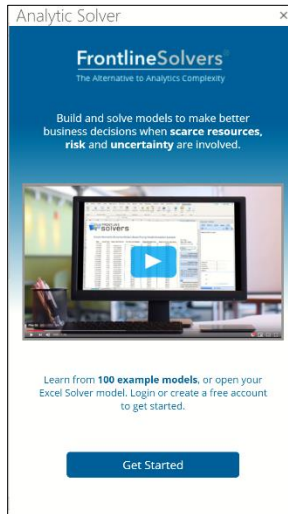
Note: Data Science Cloud will defer to Analytic Solver Desktop. If you find that Data Science Cloud does not appear after following these steps, you'll first need to remove Analytic Solver Desktop before installing the Cloud app. To do so simply open Excel and click File – Options – Add-ins – Com Add-ins – Go, select Analytic Solver Add-in, then click Remove. Once Analytic Solver Desktop is removed, then go back to Office Add-ins – My Add-ins and add in Data Science Cloud.

Single Sign On Functionality

Analytic Solver Cloud includes Single Sign On functionality which automatically logs in users to their Analytic Solver Cloud account using their Microsoft 365 credentials. This means that if you've signed in to your Microsoft 365 account using the same email address you used to register on

www.solver.com, then you will not be asked to login to Analytic Solver Cloud. Once you insert Analytic Solver cloud, you will have immediate access to all the features and functionality of Analytic Solver. As long as you remain signed in to your Microsoft 365 account, you'll never have to login to Analytic Solver Cloud again!

If you log out of Analytic Solver or your Office 365 credentials do not match your Solver credentials, you'll see the following welcome screen in the Solver task pane. Just click the Get Started button or click License – Login/Logout on the ribbon to login.

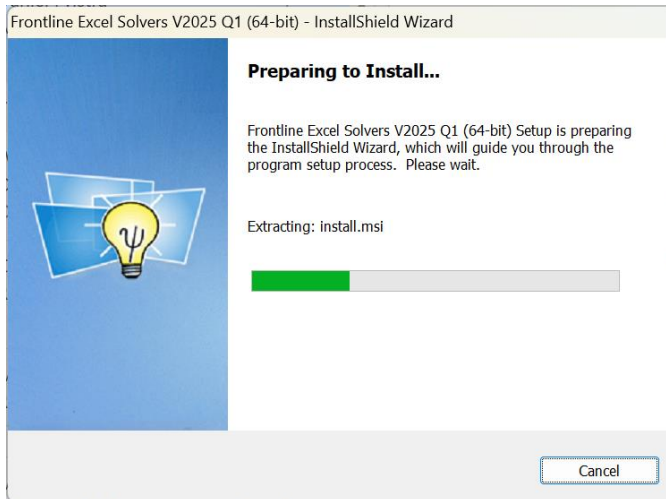


Installing Analytic Solver Desktop

To install Analytic Solver Desktop to work with any of the supported versions of Microsoft Excel (see above), simply run the program **SolverSetup.exe**, which installs all of the Solver program, Help, User Guide, and example files. SolverSetup.exe checks your system, detects what version of Office you are running (32-bit or 64-bit) and then downloads and runs the appropriate Setup program version.

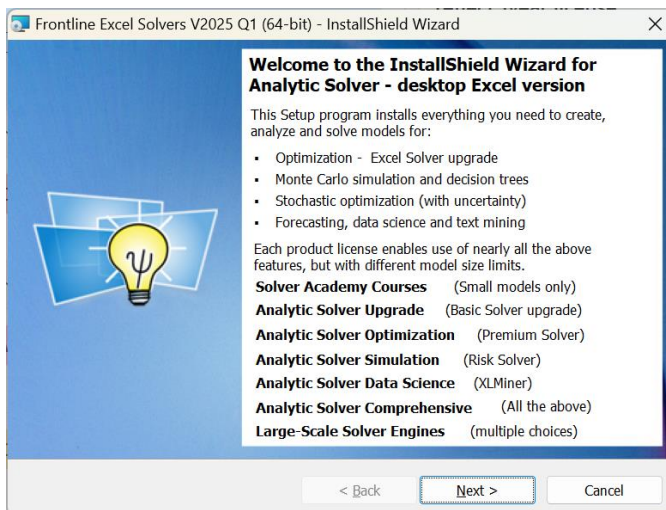
Note that your copy of the Setup program will usually have a filename such as SolverSetup_12345.exe; the '12345' is your user account number on Solver.com.

When you run the Setup program, depending on your antivirus program or Windows security settings, you might be prompted with a message “[reason such as new/unknown program]. Are you sure you want to run this software?” You may safely click **Run** in response to this message.

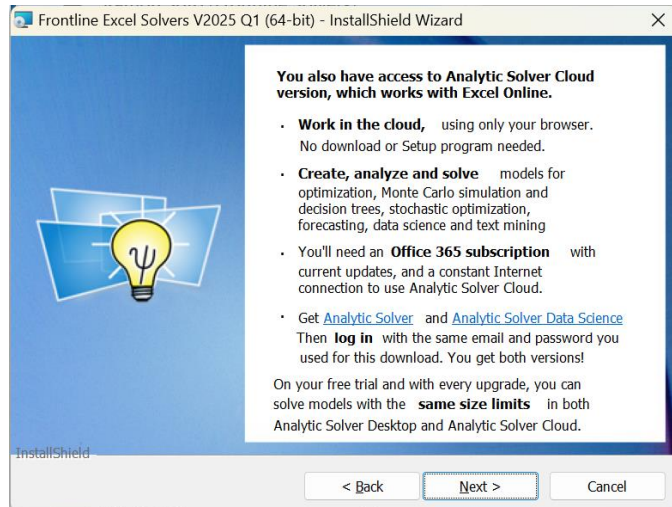


Next, you'll briefly see the standard Windows Installer dialog.

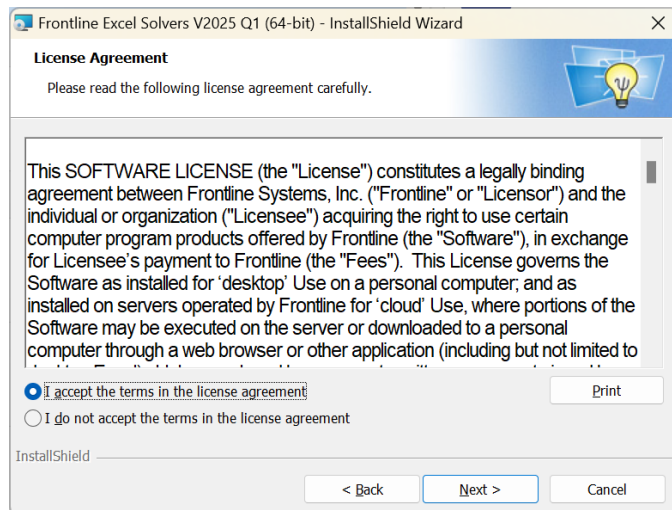
Then a dialog box like the one shown below should appear:



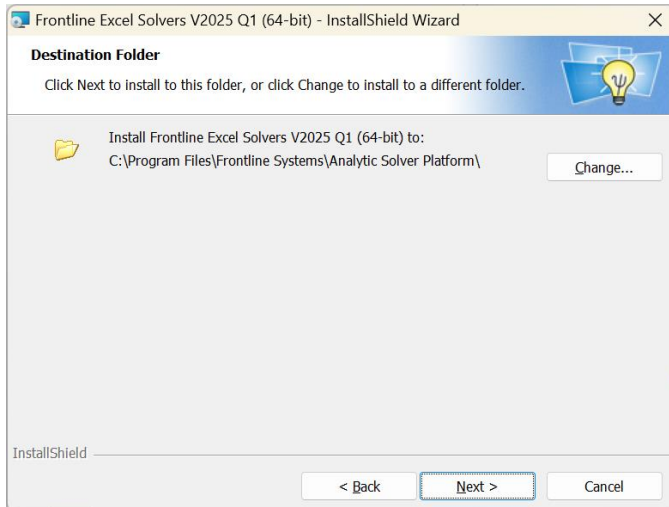
Read this, so you know the *difference* between Analytic Solver Comprehensive and its subsets. Then click **Next** to proceed – you'll see a dialog like the one below.



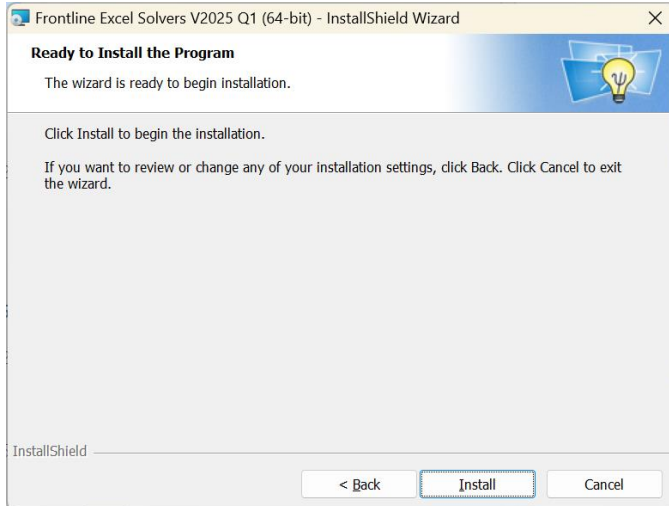
Next, the Setup program will ask if you accept Frontline’s software license agreement. You must click “I accept” and **Next** in order to be able to proceed.



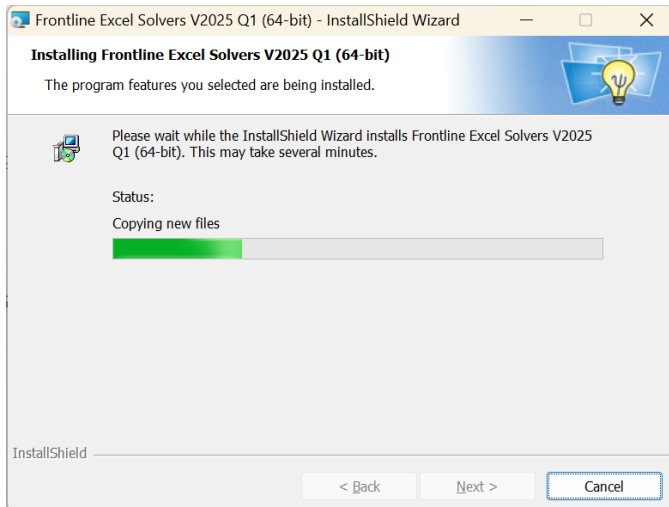
The Setup program then displays a dialog box like the one shown below, where you can select or confirm the folder to which files will be copied (normally C:\Program Files\Frontline Systems\Analytic Solver Platform, or if you’re installing Analytic Solver for 32-bit Excel on 64-bit Windows, C:\Program Files (x86)\Frontline Systems\Analytic Solver Platform). Click **Next** to proceed.



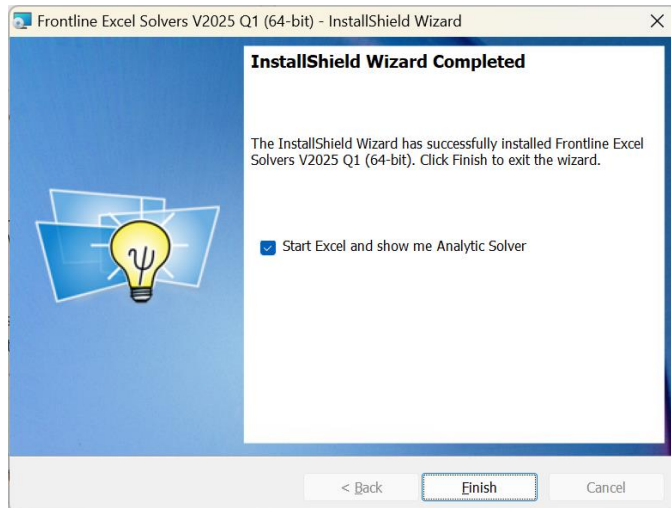
You'll see a dialog confirming that the preliminary steps are complete, and the installation is ready to begin.



After you click **Install**, the Analytic Solver files will be installed, and the program file RSPAddin.xll will be registered as a COM add-in (which may take some time). A progress dialog may appear; be patient, since this process could take longer than it has in previous Solver Platform releases.



When the installation is complete, you'll see a dialog box like the one below. Click **Finish** to exit the installation wizard.



The full Analytic Solver product family is now installed. With your trial and paid license, you can access every feature of the software, including forecasting and data science, simulation and risk analysis, and conventional and stochastic optimization. Simply click “Finish” and Microsoft Excel will launch with a Welcome workbook containing information to help you get started quickly.

1 A B C D E F G H I J K L M

2  **FRONTLINE solvers** *Developers of the Excel Solver*

3 **Welcome to your Free Trial of Analytic Solver - our tools for Optimization, Simulation/Risk Analysis, and Data Science.**

4 **Here are a few tips to help you get started:**

6 **1. Look for the *Add-Ins, Analytic Solver and Data Science* Tabs on the Ribbon above:**

7 

11 The **Analytic Solver** tab provides a new **Ribbon** for optimization, simulation, decision trees and sensitivity analysis, plus a **Task Pane** showing your optimization or simulation model. The **Data Science** tab provides a new **Ribbon** for data visualization, time series forecasting, data mining and text mining.

15 The **Add-Ins** tab contains a **Premium Solver** button which displays a **Solver Parameters** dialog similar to the basic Excel Solver we developed for Microsoft. Changes you make here are reflected in the Analytic Solver Task Pane, and vice versa.

19 **2. Start with a free trial of Analytic Solver, which includes everything in optimization, simulation, and data science.**

20 You can get *complete* model and solution information for all *sized* models. Open the License Center (License button) to purchase a license for Analytic Solver Upgrade, or our Optimization, Simulation, or Data Science full versions - pay for only what you need. **Analytic Solver Comprehensive** includes full versions of everything, and plug-in large-scale Solver Engines extend optimization to handle up to millions of variables. Click the Product Guide tab to run the **Product Selection Wizard** to help you decide which version is best for you.

27 **3. When building your model, try the Distribution Wizard (for simulation, click the Distributions button) or Constraint Wizard (optimization).**

28 

35 **4. Get started with our short guides:** [Quick Start Guide](#) [Excel Solver Upgrade Guide](#)

36 Learn to use the Ribbon and Task Pane interface, get Help, deal with licensing, and build your first model. If you're upgrading from the Solver included in Microsoft Excel, the Excel Solver Upgrade Guide can help.

39 **5. Questions?** Use the **AI Agent** or explore other resources in the **Help Center**.

Logging in the First Time

In Analytic Solver V2025 Q1, your license is associated with **you**, and may be used on **more than one PC**. For example, you can run SolverSetup to install the desktop software on your office PC, your company laptop, and your PC at home. But only **you** can use Analytic Solver, and only on **one** of these computers **at a time**. It is unlawful to “share” your license with another human user.

The first time you run Analytic Solver (Desktop or Cloud) after installing the software on a new computer, when you next start Excel and visit the Analytic Solver tab on the Ribbon, **you will be prompted to login**. Enter the **email address** and **password** that you used to register on Solver.com. Once you’ve done this in Analytic Solver Desktop, your identity will be “remembered,” so you won’t have to login every time you start Excel and go to one of the Analytic Solver tabs. In Analytic Solver Cloud, you may be asked to login more frequently.

You can login and logout at any time, by visiting the **Solver Home tab** (see “Cloud Version and Solver Home Tab” below) and clicking the **Log In** or **Log Out** button in Analytic Solver Desktop or by clicking **Help – Login/Logout** in Analytic Solver Cloud. If you share use of a single physical computer with other Analytic Solver users, be careful to **login** with your own email and password, and **log out when you’re done** – if you don’t, other users could access private files in your cloud account, or use up your allotted CPU time or storage.

When you move from one computer to another, you should **log out** on one and **log in** on the other. As a convenience, if you log in to Analytic Solver on a new computer when you haven’t logged out on the old computer, Analytic Solver will let you know, and offer to automatically log you out on the other computer.

Uninstalling the Software

To uninstall Analytic Solver Desktop, just run the **SolverSetup** program as outlined above. You'll be asked to confirm that you want to remove the software.

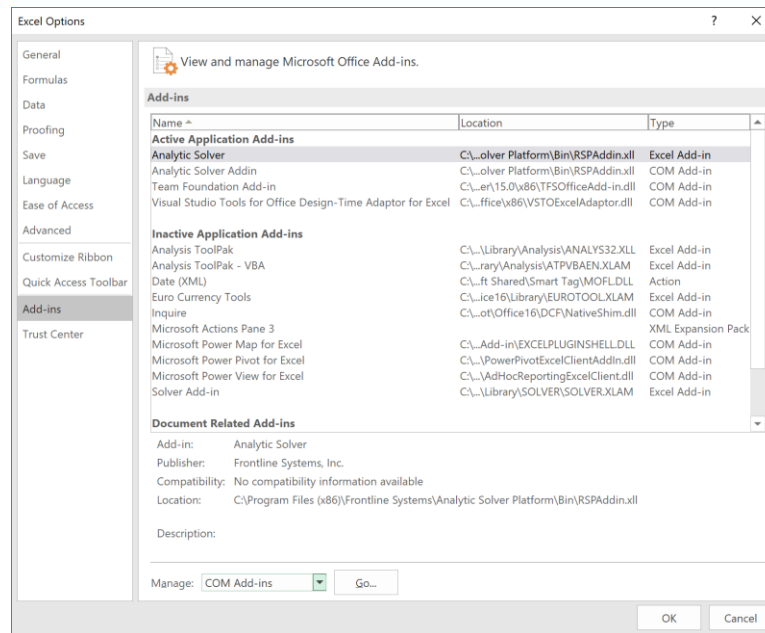
You can also uninstall by choosing **Control Panel** from the **Start** menu, and double-clicking the **Programs and Features** or **Add/Remove Programs** applet. In the list box below "Currently installed programs," scroll down if necessary until you reach the line, "Frontline Excel Solvers V2025 Q1," and click the **Uninstall/Change or Add/Remove...** button. Click **OK** in the confirming dialog box to uninstall the software.

Activating and Deactivating the Software

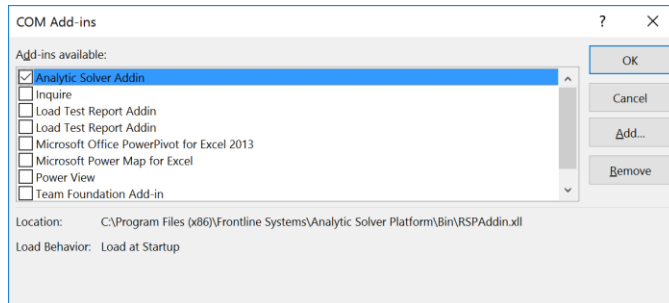
Analytic Solver Desktop's main program file **RSPAddin.xll** is a COM add-in, an XLL add-in, and a COM server. A reference to the add-in **Solver.xla** is needed if you wish to use the "traditional" VBA functions to control Analytic Solver, instead of its new VBA Object-Oriented API.

Excel 2019, 2016 and 2013

In Excel 2019, 2016 and 2013, you can manage all types of add-ins from one dialog, reached by clicking **File – Options -- Addins**.



You can manage add-ins by selecting the type of add-in from the dropdown list at the bottom of this dialog. For example, if you select **COM Add-ins** from the dropdown list and click the **Go** button, the dialog shown below appears.



If you uncheck the box next to “Analytic Solver Addin” and click OK, you will deactivate the Analytic Solver COM add-in, which will remove the Analytic Solver tab from the Ribbon in desktop Excel, and also remove the PSI functions for optimization from the Excel Function Wizard.

Analytic Solver Data Science Overview

Overview

This Guide shows you how to use **Analytic Solver Data Science Desktop** or **Cloud** (formerly XLMiner or more recently Analytic Solver Data Mining) that combines the capabilities of data analysis, time series analysis, classification techniques and prediction techniques. Analytic Solver Data Science is included in Analytic Solver Comprehensive, or can be purchased as a stand-alone license.

Analytic Solver Data Science was originally developed and marketed by others (Cytel Corp and Statistcs.com), and was in popular use for over a decade, primarily in teaching, when Frontline Systems acquired all rights to the product in 2011. In 2012-2013, Frontline marketed and supported the existing product, while rewriting the underlying software from the ground up.

- In our V2014 release, we introduced a fundamental new capability in Analytic Solver Platform, Risk Solver Platform and Premium Solver Platform for building optimization and simulation models in Excel: Dimensional Modeling. It introduces new concepts such as dimensions and cubes, and provides all the tools you need to build and solve larger scale, better structured, more maintainable models using these concepts.
- In our V2014-R2 release, Analytic Solver Platform included a completely re-engineered, far more powerful data science and forecasting capability named XLMiner Platform. New data science algorithms are up to 100 times faster, constantly exploit multiple processor cores, and offer greater accuracy and numeric stability.
- Our V2015 release introduced a wide range of new features, including powerful text mining and ensemble methods for classification and prediction in XLMiner Platform; feature selection, partitioning “on-the-fly,” ROC/RROC curves, enhanced linear and logistic regression, and more in XLMiner Pro and Platform; extensive chart enhancements, distribution fitting, and new Six Sigma functions in Risk Solver Pro and Platform; and support for “publishing” optimization and simulation models to Excel for the Web (formerly Excel Online) and Google Sheets. This feature was further extended to handle large-scale models in V2015-R2.
- Our V2015-R2 release made it easy to share results in the cloud: You can transfer optimization, simulation, or data mining results into your Microsoft Power BI online account, visualize those results with just a few clicks, and share them with others. Similarly, you can export optimization, simulation or data mining results into Tableau Data Extract (*.tde) files that can be opened in Tableau, the popular interactive data visualization software. V2015-R2 also links your Excel workbook with “Big Data”: You can easily obtain sampled and summarized results drawn from the largest datasets, stored across many hard disks, in compute clusters running Apache Spark.

V2016 Release

Our V2016 release dramatically speeds the process of moving from an “analytic Excel model” to a “deployed application, available to others.” With the new **Create App** feature, which translates your Excel optimization or simulation model into Frontline’s new **RASON** modeling language, you can create an application that can run in a **web browser**, or a **mobile app** for phones or tablets – with just two mouse clicks! Your app solves problems via our RASON server, running 24x7 on Microsoft Azure, using its REST API.

Other enhancements in V2016 include a Task Pane **navigator for data mining** in XLMiner Pro and Platform, and a greatly enhanced **Evolutionary Solver** in Premium Pro and Platform, and Risk Solver Pro and Platform. A completely new local search algorithm called **SQP-GS** (Sequential Quadratic Programming with Gradient Sampling), and new **Feasibility Pump** methods for both continuous and integer variables help the Evolutionary Solver find better solutions, faster than ever.

V2016-R2 Release

In our V2016-R2 release, simulation models can use **compound distributions**, with either a constant or a discrete distribution as the ‘frequency’ element, and **correlation using copulas** (Gaussian, Student and Archimedean forms), as well as rank-order correlation, to generate samples for multiple uncertain variables.

In optimization, the **Evolutionary Solver** is further enhanced with new ‘GA methods’ for integer variables, often yielding much better integer solutions than previous releases in a given amount of time. And users of **Dimensional Modeling** will see major improvements in speed and memory use, thanks to new support for ‘sparse cubes’. Export of analytic model results to Tableau data visualization software is more convenient than ever in V2016-R2 with support for the **Tableau Web Data Connector** introduced in Tableau 9.1.

V2016-R3 Release

With our V2016-R3 release for Excel, we’re introducing **AnalyticSolver.com**, a new cloud-based platform for both predictive and prescriptive analytics models that you can use via a **web browser** – including all the optimization, simulation, and data mining power found in our desktop products. Your models are solved using our Azure-based RASON cloud servers.

The AnalyticSolver.com **user interface** works just like our Excel user interface, with a Ribbon, Analytic Solver Platform and XLMiner Platform tabs with the same icons and dropdown menus, and a Task Pane with Model, Platform, Engine and Output tabs. Both V2016-R3 in Excel and AnalyticSolver.com include a new “Solver Home” tab on the Ribbon that makes it easy to move Excel workbooks and other files between desktop and cloud. And access to AnalyticSolver.com is **included** with your V2016-R3 license for desktop Excel.

V2017 Release

In V2017, we are introducing commercial users to **AnalyticSolver.com**, a new cloud-based platform for both predictive and prescriptive analytics models that you can use via a **web browser** – including all the optimization, simulation, and data mining power found in the desktop version. The AnalyticSolver.com **user interface** works just like our Excel user interface, with a Ribbon and Task Pane. Both V2017 in Excel and AnalyticSolver.com include a new “Solver Home” tab on the Ribbon that makes it easy to move Excel workbooks and other files between desktop and cloud. And access to AnalyticSolver.com is *included* with your V2017 license for desktop Excel.

V2017 also uses a **new licensing system** that offers you more flexible ways to use the software, both desktop and cloud. Your license is associated with **you**, and may be used on **more than one PC**. For example, you can install the software on your office PC, your company laptop, and your PC at home. But only **you** can use Analytic Solver, and only on **one** of these computers **at a time**.

V2017 introduces **Analytic Solver Basic**, as described above, to give you access to **all** Analytic Solver features, **all** the time, for learning purposes using small models. It also includes a new **License/Subscription Manager** and a **Product Selection Wizard** that makes it much easier to upgrade or change your license subscription on a self-service basis, and a new **Test Run/Summary** feature that lets you see exactly how your model will run with an Analytic Solver upgrade, even a plug-in large-scale Solver Engine, before you purchase the upgrade – and do this any time, not just during a 15-day free trial.

V2017 includes major enhancements to **data mining**: Automatic support for **categorical variables** in many classification and prediction algorithms that ‘normally’ require continuous variables; **ensembles** that combine nearly any type of algorithm as a ‘weak learner’, not just for example trees; general-purpose **Rescaling** as a new Data Transformation method that can also be applied ‘on-the-fly’ when training a model; greatly enhanced multilayer **neural networks**; ability to export models in **PMML**; and many report and chart enhancements.

The V2017 **Evolutionary Solver** includes another set of major enhancements in its handling of non-smooth models with integer variables – enough so that *most* such models will solve *significantly* faster. And there’s support for the Tableau Web Data Connector 2.0, and a new SolverSetup program that *automatically* installs the correct **32-bit or 64-bit** version of the software.

V2017-R2 Release

V2017-R2 includes major enhancements to Monte Carlo simulation/risk analysis and optimization. It’s now possible to **fit copula** parameters to historical data – a complement to distribution fitting that is sometimes called “correlation fitting.” You can use a new family of probability distributions, called the **Metalog distributions**, even more general than the Pearson distributions – members of the family can be chosen based directly on historical data (even just a few observations), without a distribution fitting process.

The V2017-R2 PSI Interpreter includes major **speed enhancements** for large linear and nonlinear **optimization models** – users with large models are likely to see a dramatic speedup in “Setting Up Problem...” Also part of this release are new, higher performance versions of the Gurobi Solver Engine (based on Gurobi

7.5), the Xpress Solver Engine (based on Xpress 30.1), and the Knitro Solver Engine (based on Knitro 10.3).

Creating Power BI Custom Visuals

The most exciting new feature of V2017-R2 is the ability to turn your Excel-based optimization or simulation model into a **Microsoft Power BI Custom Visual**, with just a few mouse clicks! Where others must learn JavaScript (or TypeScript) programming and a whole set of Web development tools to even begin to create a Custom Visual, you'll be able to create one right away.

You simply select rows or columns of data to serve as changeable parameters, then choose **Create App – Power BI**, and save the file created by V2017-R2. You click the Load Custom Visual icon in Power BI, and select the file you just saved. What you get isn't just a chart – it's your *full optimization or simulation model*, ready to accept Power BI data, **run on demand** on the web, and display visual results in Power BI! You simply need to drag and drop appropriate Power BI datasets into the “well” of inputs to match your model parameters.

How does *that* work? The secret is that V2017-R2 translates your Excel model into **RASON®** (RESTful Analytic Solver Object Notation, embedded in JSON), then “wraps” a JavaScript-based Custom Visual around the RASON model. See the chapter “Creating Your Own Application” for full details!

V2018 Release

V2018 extends Analytic Solver's forecasting and data mining features with a new capability called **data mining workflows** that can save a lot of time and eliminate repetitive steps. You can combine nearly any of Analytic Solver's data retrieval, data transformation, forecasting and data mining methods into a single, all-inclusive workflow, or pipeline.

Using the new Workflow tab in the Task Pane, you can either “**drag and drop**” icons onto a “canvas” to create a workflow diagram, or you can simply turn on a **workflow recorder**, carry out the steps as you've always done by choosing menu options and dialog selections, and the workflow diagram will be created automatically. Once the diagram or pipeline is created, you can “run” it in one step – each data mining method in the workflow will be executed in sequence.

In previous releases, you could use the trained model from a *single* data mining method (such as a Classification Tree or Neural Network) to “score” new data, by mapping features (columns) between the training set and new data set. In V2018, you can apply an *entire workflow* – including data transformations, partitioning, model training, and more – to a new dataset, by mapping features (columns) between the dataset used to create the workflow, and a new dataset.

Creating Tableau Dashboard Extensions

Another exciting new feature of V2018 is the ability to turn your Excel-based optimization or simulation model into a **Tableau Dashboard Extension**, with just a few mouse clicks! This is quite similar to the ability to create Power BI Custom Visuals introduced in Analytic Solver V2017-R2. It works (only) with Tableau version 2018.2 or later.

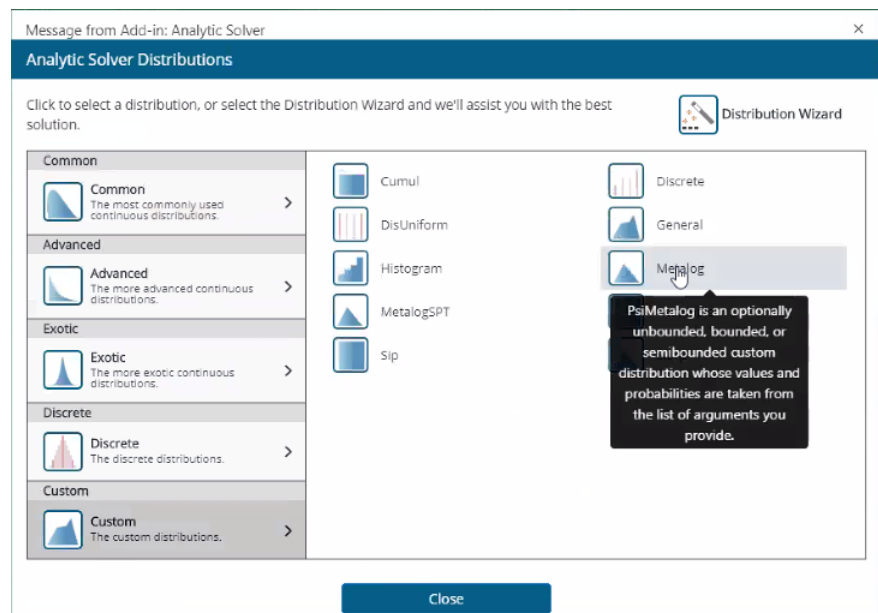
You simply select rows or columns of data to serve as changeable parameters, then choose **Create App – Tableau**, and save the file created by V2018. In Tableau, drag the **Extensions** object onto your dashboard, and choose the file you just saved. You'll be prompted to match the parameters your model needs with data in Tableau. What you get isn't just a chart – it's your *full optimization*

or simulation model, ready to accept Tableau data, **run on demand** (using our **RASON** server), and display visual results in Tableau!

V2019 Release

In Frontline Solvers V2019, we are introducing **Analytic Solver Cloud** – a next-generation product that’s the result of five years of development using new cloud technologies, that we can now bring to you – since Microsoft has released a complete set of JavaScript APIs for new Excel features, such as functions like PsiNormal() and PsiMean() used in simulation and risk analysis.

While we have new and enhanced features in development, in the V2019 release we have focused on a **consistent user experience** between Analytic Solver Desktop and Analytic Solver Cloud. In a few cases, this has involved modifications to Analytic Solver Desktop. For example, older versions of Analytic Solver Desktop used “cascading submenus” to select probability distributions, and results in Monte Carlo simulations. Since a JavaScript add-in cannot define or use cascading submenus, we have modified Analytic Solver Desktop *and* Analytic Solver Cloud so that the **Distributions** button on the Ribbon displays a dialog, rather than a cascading submenu, where you can select an appropriate probability distribution. Yet the order and layout of probability distributions remains the *same* as in previous Analytic Solver versions.



V2020 Release

With each new release, Analytic Solver gives you more! Our latest enhancements apply to both Analytic Solver Desktop and Analytic Solver Cloud.

In V2020, the LP/Quadratic Solver – probably the most-used Solver Engine in Analytic Solver – features significantly improved performance on linear mixed-integer models. Prior versions of this Solver would use only one processor core at a time, but V2020 will use *multiple processor cores* to speed your solution.

The plug-in large-scale Solver Engines in Analytic Solver V2020 also feature significantly improved performance (they continue to utilize multiple processor

cores). These include the **Gurobi** Solver V9.0, with a new ability to solve non-convex quadratic models; the **Xpress** Solver V35, with a new Solution Refiner, and the Artelys **Knitro** Solver V12.1, with SOCP and MIP speedups.

Analytic Solver V2020 includes 12 new probability distribution functions, *enhanced* property functions for the PSI Distribution functions, *new* property functions for the PSI Statistics functions, and *new* “theoretical” functions that return analytic moments of distributions. Full details are in the *Frontline Solvers Reference Guide*, but here’s a partial list of new/enhanced functions:

PsiBurr	PsiLevyAlt	PsiTheoMin
PsiDagum	PsiHypSecantAlt	PsiTheoMax
PsiDbfTriang	PsiCumulD	PsiTheoVariance
PsiFatigue	PsiLaplace	PsiTheoStdDev
PsiFdist	PsiCauchy	PsiTheoSkewness
PsiFrechet	PsiTruncate	PsiTheoKurtosis
PsiHypSecant	PsiCensor	PsiTheoRange
PsiJohnsonSB	PsiLock	PsiTheoPercentile
PsiJohnsonUB	PsiOutput	PsiTheoPercentileD
PsiKumaraswamy	PsiTheoMean	PsiTheoTarget
PsiReciprocal	PsiTheoMedian	PsiTheoTargetD
PsiLevy	PsiTheoMode	PsiCategory

These functions make it **even easier** to adapt risk analysis models developed with other popular Excel add-ins to **work with Analytic Solver**. A simple Find and Replace of the function name prefix with ‘Psi’ is often all you need. And unlike those other Excel add-ins, with Analytic Solver you can easily run your model in the cloud with **Excel for the Web**, translate your model to RASON, and use it in **Power BI, Tableau**, or your own **web or mobile** application!

What’s New in Analytic Solver V2020.5

Analytic Solver V2020.5 includes significant enhancements to both Monte Carlo simulation and optimization – but the most exciting new feature is a greatly expanded **Create App** facility that makes it easy to **deploy** your Excel analytic model as a **cloud service** (thanks to RASON), usable from nearly any corporate, web or mobile application. What’s more, you can manage, monitor and update your own cloud services, without ever leaving Excel!

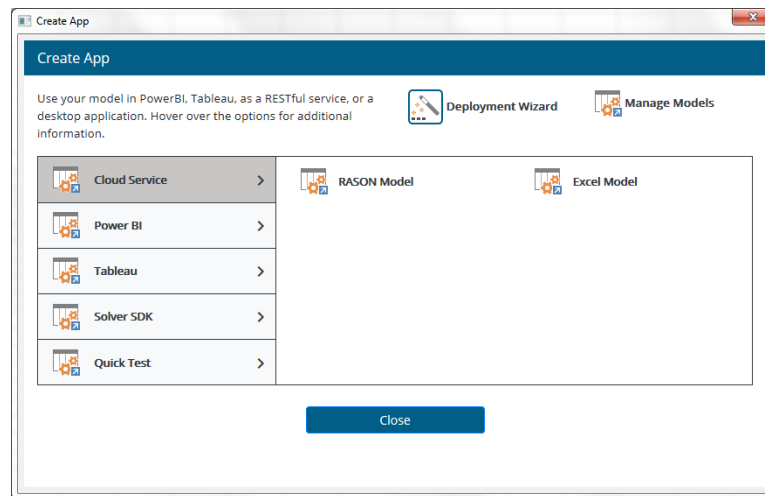
Easily Deploy Your Model as a Cloud Service

We’ve realized for many years that developing and testing your analytic model in Excel is often just the first step: To gain the **real business value** from the decisions that can be made with your model, it’s often necessary to get the model into the hands of **other people in the business** – in a form where they can easily ensure that it has **up-to-date data, re-run** the model’s optimization, simulation, or data mining process, and either **view** the results, or **plug them into** another software application or process.

In our 2017 and 2018 Analytic Solver releases, we took the steps that were possible at that time, enabling users to get their models into the hands of **Power**

BI and **Tableau** users. And we built a facility to *translate* simpler Excel models into the RASON modeling language, enabling them to be solved in our cloud platform (RASON is an acronym for RESTful Analytic Solver Object Notation). But up to this point, a typical Excel user would still need help from a web developer, or would need learn JavaScript and other web development skills, to make truly effective use of this facility.

Now in our V2020.5 release of both Analytic Solver and RASON, we've gone much further to **simplify** the process of **deploying** an Excel model as a cloud service, and **connecting** it to databases and cloud data sources. The RASON cloud service will now accept and run **Excel workbook** models “on a par” with models written in the RASON modeling language. With the **Create App** menu option, you can turn your Excel model into a cloud service in seconds.



As an Analytic Solver user, you can now create and test models, deploy them “to the cloud” – **point and click** – as full-fledged RESTful decision services, and even get reports of recent runs of your decision services, all without leaving Excel. Using our web portal at <https://rason.com>, you can go further –even embed your Excel workbook in a multi-stage “decision flow” that can combine SQL, RASON, Excel, and DMN models, passing results from stage to stage.

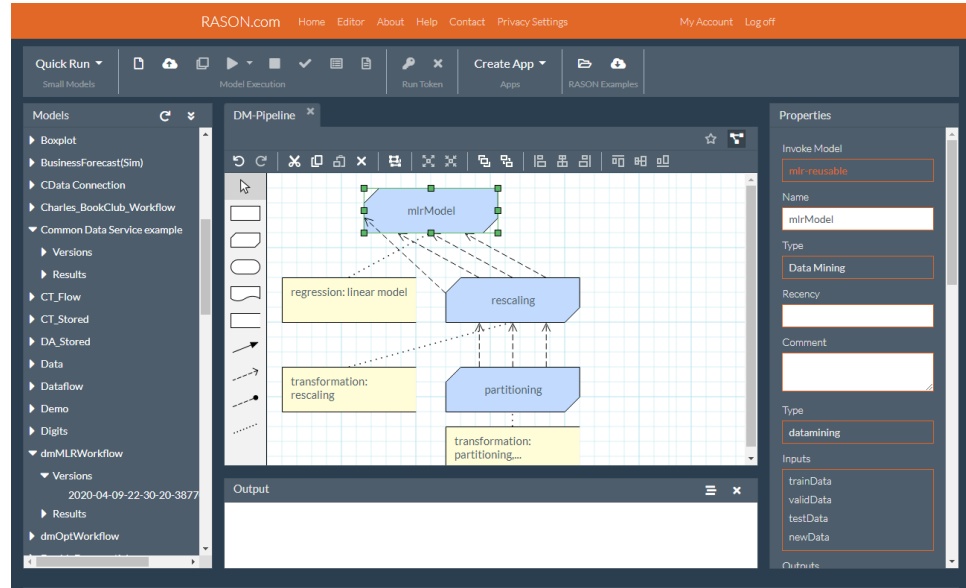
We start you out with a **RASON trial license**, so you can try out these new capabilities without purchasing anything else! (Of course, you may need an upgraded RASON license to deploy your model to many users, and re-solve it hundreds or thousands of times on our cloud servers.)

More About RASON Decision Services

RASON is an Azure-hosted cloud service that enables your company to easily embed 'intelligent decisions' in a custom application, manual or automated business process, applying the full range of analytics methods – from simple calculations and business rules to data mining and machine learning, simulation and risk analysis, and conventional and risk-based optimization.

RASON Decision Services can be used from nearly **any application**, via a series of simple REST API requests to <https://rason.net>. To express the full range of analytic models, RASON includes a high-level, declarative modeling language, syntactically embedded in JSON (JavaScript Object Notation), the popular structured format almost universally used in web and mobile apps. RASON results appear in JSON, or as more structured OData JSON endpoints.

RASON Decision Services also includes comprehensive data access support for **Excel, SQL Server on Azure, Power BI, Power Apps, Power Automate** (aka Microsoft Flow) and **Dynamics 365**. And it includes powerful model management tools, such as tracking model versions including “champions and challengers”, monitoring model results, and automated scheduling of runs for both models and multi-stage decision flows.



How You Can Use RASON

You can use RASON to quickly and easily create and solve optimization, simulation/risk analysis, data mining, decision table, and decision flow models – instantly deployed as cloud services. You can learn RASON, create models, supply data and solve them, and even manage model versions and cloud data connections, “point and click” using <https://rason.com>, our “web portal” to the underlying REST API service.

If you’ve used another **modeling language** to build an analytic model, you’ll find the RASON language to be simple but powerful and expressive – and integrating RASON models into a larger application, especially a web or mobile app, is **much easier** than with other modeling languages. Excel users will find that RASON includes virtually the entire **Excel formula language** as a subset. If you’ve used tools based on the DMN (Decision Model and Notation) standard, you’ll find that RASON – and Analytic Solver, as shown in the chapter “Building Decision Tables” – fully support DMN and FEEL Level 2.

Unlike existing “heavyweight” Business Rule Management Systems, with year-long implementation schedules, six-figure budgets and limited analytics power, RASON Decision Services enables you to **get results in just weeks to months**, from building and testing models, to deploying them across an organization. With RASON, you can build successful POCs (Proofs of Concept) without any IT or professional developer support – yet RASON is very “IT and developer friendly” when you’re ready to deploy your POC across your company.

New Time Series Simulation Functions

Analytic Solver V2020.5 includes another new set of PSI Distribution functions and related PSI property functions, focused around **time series simulation**. Earlier Analytic Solver versions supported time series simulation using functions such as PsiForecast() and PsiPredict(), and models fitted via Analytic Solver Data Mining – but V2020.5 goes further, to support time series functions found in other popular Excel add-ins, such as Palisade’s @RISK. Full details are in the *Frontline Solvers Reference Guide*, but here’s a partial list of new/enhanced functions:

PsiAR1	PsiBMMR	PsiAPARCH11
PsiAR2	PsiGBMJD	PsiTSTransform
PsiMA1	PsiARCH1	PsiTSIntegrate
PsiMA2	PsiGARCH11	PsiTSSeasonality
PsiARMA11	PsiEGARCH11	PsiTSSync

With these functions, virtually any risk analysis model developed with other popular Excel add-ins, such as Palisade’s @RISK, can be easily made to work with Analytic Solver. An appendix in the *Frontline Solvers Reference Guide*, “@Risk to Analytic Solver Psi Function Conversion Table”, explains the details. And with Analytic Solver, you can easily **deploy** your risk analysis model as a **cloud service** – usable from Tableau, Power BI, Power Apps, Power Automate, or virtually any corporate, web or mobile application!

New Optimization Result Functions

In every version of Analytic Solver (and its predecessors, such as Premium Solver and the Excel Solver), you could obtain all the properties of an optimal solution – such as initial and final values, dual values, and ranges for decision variables and constraints – via the Answer Report and Sensitivity Report, which are inserted into your Excel workbook as new worksheets. But what if you want only **select subsets** of these values – and you’d like to have them on the **same worksheet** as your model? That’s now possible in Analytic Solver V2020.5. Just type these new functions into cells, or use the Function Wizard in Excel.

PsiInitialValue	PsiDualValue	PsiCalcValue
PsiFinalValue	PsiDualLower	PsiOptStatus
PsiSlackValue	PsiDualUpper	PsiModelDesc

These new functions have another purpose in V2020.5, when you use Analytic Solver’s enhanced **Create App** facility to deploy your model as a cloud service: You can use them to determine a select subset of values from the solution that you want to **return from your cloud service** to a calling application.

What’s New in Analytic Solver V2021

Analytic Solver V2021, our latest release, features faster startup and better support for 4K monitors in Analytic Solver Desktop, improved support for decision trees "on the spreadsheet" in Analytic Solver Cloud, and a new, faster version of the Xpress Solver (V37.1.3) – but the most exciting new features in this release give you new ways to define your own **custom functions**, in a way that works in both Excel Desktop and Excel for the Web, as well as our cloud platform RASON.

Lambda, Let and Box Functions

In years past, Excel-savvy analysts used VBA (Visual Basic for Applications) to define their own custom functions. While this still works in Excel Desktop using COM (28-year-old Component Object Model), VBA functions are *not* supported in Excel for the Web – and according to Microsoft, VBA and COM will *never* move to the cloud. If you want your custom functions to work in *both* desktop and cloud, your options have been limited – until now:

- Microsoft has introduced new Excel functions LAMBDA and LET. These are very special because you can use them in Excel formulas to define your own custom functions. The Excel community has expressed much excitement over these new functions, since they effectively make Excel a “complete programming language”. (In Q1 2021, these functions are being rolled out across the different Office update channels.)
- On another front, there’s the open standard known as DMN (Decision Model and Notation) – a business user-friendly “formula language” used to define business rules and decision tables, supported in “decision management” platforms from various vendors, and in Analytic Solver and RASON since 2019. DMN – now in version 1.3 – offers a way to define your own custom functions, known as “Box functions”.

Analytic Solver V2021 includes support for *both* Excel’s LAMBDA and LET functions, and for DMN-compatible Box functions. You’ll find a new chapter in this User Guide, “Using Custom Functions”, that explains how to use both approaches. Even better, LAMBDA and LET, and DMN Box functions enjoy full support from our PSI Interpreter – which means that our full range of Solver Engines, and our high-speed Monte Carlo simulation engine “understand” and take full advantage of custom functions that you define this way. This can yield better results than you’ve ever had with VBA-based functions that are embedded in an optimization or simulation model.

What’s New in Analytic Solver V2021.5

Automate Data Mining with Find Best Model

Analytic Solver, our SDKs and RASON include comprehensive, powerful support for data mining and machine learning. Using these tools, you can “train” or fit your data to a wide range of statistical and machine learning models: Classification and regression trees, neural networks, linear and logistic regression, discriminant analysis, naïve Bayes, k-nearest neighbors and more. But the task of choosing and comparing these models, and selecting parameters for each one was up to you.

With the new Find Best Model options in V2021.5, you can automate this work as well! Find Best Model uses methods similar to those in (expensive high-end) tools like DataRobot and RapidMiner, to automatically choose types of ML models and their parameters, validate and compare them according to criteria that you choose, and deliver the model that best fits your data.

Better Simulation Models with Metalog Distributions and Fitting

Analytic Solver, our SDKs and RASON support over 60 “classical” probability distributions for Monte Carlo simulation. Since mid-2017, they’ve also supported the increasingly popular Metalog family of distributions, created by

Dr. Tom Keelin, and recently popularized by the nonprofit Probability Management group. Metalog distributions can closely approximate virtually any classical continuous distribution, and often they can better fit user data than classical distributions. In V2021.5 we've brought Metalog distributions to the fore, with a powerful new facility to automatically fit user data to the full range of possible (bounded and unbounded, multi-term) Metalog distributions. It's never been easier to get an accurate probability distribution that fits a real-world phenomenon.

Share Data Mining and Probability Models via RASON

In recent years, our Azure-hosted RASON Decision Services platform has offered increasingly powerful facilities to deploy models to the cloud and share them with other users – culminating in Frontline Solvers V2020.5, when we enabled deployment, sharing, versioning and management of Excel models as well as native RASON models, plus support for multi-stage decision flows, encompassing and going beyond traditional data science workflows.

In V2021.5 we've gone further: You can now deploy and share data mining and machine learning models, trained in Analytic Solver or RASON, to the Azure cloud, and use them directly for classification and prediction (without needing auxiliary “code” in R or Python, RASON or Excel). You can also deploy and share probability models, following the open Probability Management 3.0 standard. Using these Shared Information Probability resources (SIPs, also known as “Stochastic Information Packets”), you can ensure that your group or organization uses consistent data about uncertain/risky variables across simulation or decision models, enabling model results to be meaningfully compared.

And of Course, Optimization Enhancements

Analytic Solver, our SDKs and RASON have always offered rich support for conventional and stochastic optimization, improved in every new release.

V2021.5 is no exception: We've made PSI Interpreter enhancements to better utilize main memory in large optimization models (with 1 million or more

decision variables) in all three product lines. In our V2021.5 release, we're also shipping the latest Gurobi Solver 9.1, Xpress Solver 37.1.3.0, and new KNITRO Solver V12.4 with each of our products.

Analytic Solver Upgrade

Analytic Solver is our comprehensive tool for predictive, prescriptive and decision analytics. In recent years we've added deep support for business rules and decision tables, following the DMN standard, to Analytic Solver, and similar facilities in RASON. Since we've matched some entire products offered by competitors in the “rules” or “decision management” world, we've brought together our decision trees, DMN business rules and decision tables, and DMN Box functions, and combined them with Premium Solver Upgrade (for optimization models) to create a new product, Analytic Solver Upgrade.

In Excel, you'll see a new icon group on the Analytic Solver Ribbon Decision Model, next to the Optimization Model and Simulation Model groups. You'll also see the Tools group options moved to the Task Pane (already true in Analytic Solver Cloud version). Analytic Solver Upgrade will still be included in Analytic Solver Comprehensive and in a full RASON license, but its DMN-related and optimization features will have size limits in Analytic Solver Simulation and Analytic Solver Data Mining.

Analytic Solver Academy – Replaces Analytic Solver Basic

We've found that many purchasers of Analytic Solver Basic, who lacked a previous analytics background and didn't also purchase learning aids such as our Solver.Academy courses, experienced limited success and often didn't renew their licenses. In V2021.5 we are replacing Analytic Solver Basic with a bundle of all four current Solver.Academy courses plus a license for Analytic Solver with basic limits, if no other Solver license has been purchased.

What's New in Analytic Solver V2022

Faster Interaction, Faster Solves

In Analytic Solver V2022, our latest release, Excel startup is **faster**, and your 'regular interaction' with Excel is **faster** when Analytic Solver Desktop is loaded. Dialog rendering is improved for users with very-high-resolution monitors, as well as monitors with limited vertical depth.

Most impactful, the process of interpreting the model – when the model type is diagnosed and when “Setting Up Problem...” appears in the Task Pane status bar – is **faster** for most models, and *significantly faster* for larger models, especially on models with very deep “chains” of formulas that depend on other formulas. We're also including new, **faster** versions of the Gurobi Solver for linear mixed-integer models, and the Knitro Solver for nonlinear models.

V2022 re-introduces the **Freeze** and **Thaw** options (now located on the Task Pane Tools tab), which allow you to share Analytic Solver models containing PSI function calls with users who don't have Analytic Solver installed. (“Freeze” will save PSI function call formulas in cell comments, and “Thaw” will restore them later as formulas.)

Deploy Model functionality in V2022 is enhanced on our RASON cloud server, and **Microsoft Teams** messaging is improved for several account types.

In case you didn't know: Recent Analytic Solver releases, including V2022, allow you to **Test Run** models that *exceed* the limits of your current license. (You don't get full results for all variables and constraints, but you do get the final objective value and solution time.) You can even use optional large-scale Solver Engines, like the Gurobi and Knitro Solvers, in a “Test Run”. To make it clear what your current license does and doesn't include, in V2022 you'll see “(Test Run)” next to the names of optional Solver Engines in the Task Pane.

What's New in Analytic Solver V2023

Automated Risk Analysis of Machine Learning Models

Besides our “usual enhancements” such as faster Solver Engines and new PSI functions like PsiCalcParam(), Analytic Solver V2023, our latest release, features an innovative (and patent pending) new capability for **automated risk analysis** of **machine learning** models on the Data Mining tab.

A further benefit of this new feature is a general-purpose, easy to use facility for **synthetic data generation**, to augment the data you already have. Analytic

Solver Data Mining and Comprehensive users are able to use these new features with size limits constrained only by memory, but *all* Analytic Solver users have access to these features with “Basic size limits” for your datasets.

Until Analytic Solver V2023, data science and machine learning (DSML) tools – including ours – had no facility for **risk analysis** of machine learning (ML) models, prior to their production use. Most tools (including ours) had facilities for ‘training’ the model on one set of data, ‘validating’ its performance on another set of data, and ‘testing’ it versus other ML models on a third set of data. But this is not *risk analysis*.

Synthetic Data Generation ‘For Free’

Synthetic data generation has come into use in recent years to augment available datasets, when the available data is limited, or is restricted by law or regulation, such as with personal health information (PHI). In Analytic Solver V2023, you have a powerful, general-purpose tool for synthetic data generation, using the new **Generate Data** button on the Data Mining Ribbon tab. So far, we’re keeping Analytic Solver current with the “state of the art”.

But with the new **Simulation** tab that you’ll find in every ML model training dialog in Analytic Solver, we’ve gone *beyond* the “state of the art”, to bring you a new way to assess your trained ML model’s performance – not to determine how well it **has performed** on data you have, but to quantify how it **may perform differently** on data it will encounter in the future. We generate synthetic data “on the fly” and use it in a Monte Carlo simulation of your ML model’s performance – and we highlight **differences** in model performance in training versus simulated production use.

The beauty of this approach is that **you don’t have to do any work** to obtain a risk analysis of your model’s performance, beyond a few mouse clicks on the **Simulation** tab to enable the analysis, which is entirely **automated**. You don’t even have to be familiar with the features of Analytic Solver Simulation to use this capability for machine learning!

(You *can* adjust several settings that involve distribution fitting, correlation and copulas, Monte Carlo sample generation, etc. – but the default settings work very well.) Typically, the analysis adds only **seconds** to a perhaps a **minute** to the time taken to train and validate your ML model. So you can make this a routine **part of the process** of training and assessing new ML models.

As noted earlier (opposite the title page) in this User Guide, we have a patent application pending titled “Automated Risk Analysis of Machine Learning Models”. But as an Analytic Solver (or Solver SDK or RASON) user, you gain first access to this new capability at no extra cost.

What’s New in Analytic Solver V2023 Q1

We released Analytic Solver V2023 in September 2022, with the innovative enhancements described above – but we surprised even ourselves by how much more we could offer our customers by December 2022! So we’ve christened this new release “Analytic Solver V2023 Q1”.

Faster LP/Quadratic Solver and Large-Scale LP/QP Solver Engine

Our “headline feature” for this release is a new, higher performance version of the built-in LP/Quadratic Solver, as well as its “cousin” the Large-Scale LP/QP

Solver Engine which has much higher problem size limits. We expect most users will see **faster – sometimes much faster** solutions with this new version, both for LP and QP (linear programming and quadratic programming) models, and for LP and QP models with integer constraints.

And there's more: Now you can solve models with (convex) **quadratic constraints**, with or without a quadratic objective, using this Solver Engine.

(There are always tradeoffs: We've found that for about 10% of models we've tested, the new version is not faster, and can even be *slower* for models with integer constraints. If you encounter this, just select the **Classic Search** option to get exactly the same performance you were getting before.)

More Plug-in Solver Engine Improvements

V2023 Q1 also includes new versions of other large-scale Solver Engines. The **Gurobi Solver Engine** is upgraded (to their V10.0) with a range of enhancements, yielding solution times faster by 3-10% to 25% on a wide range of linear and quadratic mixed-integer models. The **Xpress Solver Engine** is upgraded (to their V9.0) with enhancements to strong branching, separation of cutting planes, and a new heuristic method, run at the branch & bound root. The **KNITRO Solver Engine** is upgraded (to their V13.2) with a range of enhancements for smooth nonlinear mixed-integer models, including new presolve, cut selection and heuristic methods.

Greatly Improved “Deploy Your Model to Teams” Capability

We learned from surveys that a large majority of our customers work in companies using **Microsoft Teams** – so we've significantly enhanced a feature introduced in Analytic Solver V2022, that makes it easy to **share** your Excel model results with colleagues in your company, using Teams. The V2022 feature was designed to share your (entire) model, but we realize that many users want or need to share just the **model results** – not the full model with all its optimization and/or simulation model elements.

In V2023 Q1, when you choose **Deploy Model** from the Ribbon and click **Teams – Teams Report**, Analytic Solver will automatically create a new workbook holding only **model results**, with external links to your **model** workbook. You can choose exactly which optimization and/or simulation results you want to include in this workbook. The new workbook will be saved online, and made available to the users you want, through a “Teams channel” that you select. Your colleagues, using **just Teams**, will be able to open the workbook and view, copy or work with the results you're providing. And perhaps the best part: When you **re-run** your optimization or simulation model with new data, the workbook in Teams will be **automatically updated** (via those external links) with the latest model results! See the chapter “Deploying Your Model” in this Guide for full details.

Risk Analysis of Machine Learning Models Created in Other Software

Analytic Solver V2023 introduced an innovative (and patent pending) new capability for **automated risk analysis** of **machine learning** models – see “What's New in Analytic Solver V2023 Q1” above for a complete summary, and see our Data Mining User Guide for full details. But we realize that many

people create and test machine learning models using other software. Those folks just don't have the ability to quantify how their ML models **may perform differently** on data they will encounter in the future ... until now. But **you can help them**, using Analytic Solver V202 Q1 (or using RASON V2023 Q1, if they prefer to use a cloud platform).

In V2023 Q1, we've made it easy to perform **automated risk analysis** of models **created in other software** and saved in **PMML** (Predictive Modeling Markup Language) format. PMML is an open standard that is widely supported by software for machine learning. Analytic Solver and RASON will also save a trained machine learning model in PMML form. But now you can bring a PMML model into Analytic Solver, plus some of the data you used to train the model (just copy the PMML text and the data onto worksheets) – then use simple menu options to quickly get insights into the future performance of this ML model, from our **automated risk analysis** methods.

There are plenty of other small enhancements and fixes in V2023 Q1 – and we have more new features “in the works”. But you can see why we felt we were ready to deliver another major version to our customers – in time for Christmas!

What's New in Analytic Solver V2023 Q3

In Analytic Solver V2023 Q3, we've made it easier than ever to create analytic models, and easier to deploy them for ongoing use, through two major new features: our new conversational “AI Agent”, with a button next to “Help” on the Analytic Solver Ribbon, and a new “Identify Inputs” feature that simplifies the task of updating key input parameters when it's time to re-run your model.

AI Agent: Ask for Help from ChatGPT “Trained on Analytic Solver”

By now, nearly everyone who's been following developments in software is aware of ChatGPT, the conversational agent using “Generative AI” methods, developed by the OpenAI nonprofit closely affiliated with Microsoft. You can ask ChatGPT questions on almost any topic – including analytic methods – and get meaningful and interesting (though not always 100% correct!) answers.

That's great – but what if you could have a ChatGPT “technical support agent” that had studied all 2,300 pages of Frontline's User Guides (including this one), Reference Guides and QuickStart Guides – everything about optimization, Monte Carlo simulation, data science and machine learning, and business rules – and was instructed to use its knowledge to answer your questions? Well, that's what you have in Analytic Solver V2023 Q3!

Actually making this work involves a fair amount of software engineering: Generative AI tools need the right “context” to respond to your question – so we've built an online resource where all those 2,300 pages of Guides are represented via “vector embeddings” that enable searches by meaning, not just “keyword matches”. Our AI Agent automatically searches this online resource to create “context” for your queries to ChatGPT – and you can also search this online resource directly.

We're using the “real” online ChatGPT 3.5 Turbo version for its full conversational capabilities, which costs us money for every query – so you will

find some limits placed on query length and number of queries per month, depending on your software license type. But we expect you'll be able to use our AI Agent to amplify your own efforts, ultimately building more capable and effective analytic models. We look forward to your feedback!

Identify Inputs: Easily Set Up your Model for Data Updates and New Solves

When you're first building and solving a model in Analytic Solver, usually your focus is on getting a solution for a specific instance of a business problem, with data you have today – gathered from one or several external sources into Excel. After some effort, you're now getting solutions from optimization, simulation, or data science and machine learning. That's a success by itself ... but it leads to a desire to “do it again (and again) in the future.”

But ... those data and parameter values that *should be updated* for a future run may be scattered around your Excel spreadsheet, which also includes cells with calculated formulas, cells with constant data that doesn't need updating, and cells with data that you want to be visible, but that doesn't actually affect or “participate in” the analytic model. How do you find *just* the parameter values that *should* be updated – and then, how do you *actually update* them?

Now you can choose **Tools – Identify Inputs** and get help doing exactly that. This tool uses our PSI Interpreter to scan your model formulas, identify and list only those cell ranges that are candidates for data updates that will affect the model when solved. You can then pick and choose from a list of cell ranges, which ones you actually need to update for a new run.

Then you can either (i) let Analytic Solver “highlight” those cells with colors and backgrounds, or even better, (ii) automatically add calls to our PsiInput() function that reference those cells. These cells will automatically appear in the Task Pane as part of your model (under the **Input Data** heading) – but when you deploy your model for use *beyond Excel*, via our RASON and Solver SDK tools, you'll have easy ways to supply new values for exactly those input cells.

If you've planned in advance for data updates and re-solves, Tools – Identify Inputs can help you check and validate your work. But if – like most of us – you were focused on getting a first-time solution and “put off” the task of data updates and re-solves, Tools – Identify Inputs can be a huge time-saver!

What's New in Analytic Solver V2024 Q2

In **Analytic Solver V2024 Q2**, our overall theme is "Core Performance Improvements" -- indeed we think every user with a model of nontrivial size and complexity should **upgrade** (at no extra cost), since you're likely to see performance improvements -- especially faster "**Setting Up Problem...**".

Solving Outside Excel

The most visible change -- though *not* the one impacting the most users -- is on the menu for the **Optimize** button on the Ribbon: the last dropdown choice, **Solve on Solver Server**, has been greatly extended. In past Analytic Solver releases, this choice allowed you to solve your model (actually perform

the optimization) **outside** the running Excel program, in a separate program called Solver Server (part of our Solver SDK product), that can run on your PC or on another PC on your local area network (LAN). But the Excel-based Analytic Solver add-in continued to "monitor" the solving process, which meant that you couldn't do *something else* with Excel, while waiting for a "long solve" to complete.

In Analytic Solver V2024 Q2, this choice allows you to solve your model (again, actually perform the optimization) **outside** of Excel, on any of three alternatives: the SDK-based **Solver Server**, our public, Azure-based **RASON server**, or on a new offering: our **Containerized RASON Server** -- which brings all the capabilities of RASON to your own PC, another physical or virtual PC on your LAN, or your own cloud account (we'll say more about this in an upcoming blog post). And unlike in previous releases, V2024 Q2 is able to "monitor" the solving process **without** blocking your use of Excel for *something else*, even working on a different workbook. When the final solution for your long-running model is available, you'll be notified in Excel, and your workbook containing the original model will be automatically updated, without disturbing anything else.

New Versions of Gurobi and OptQuest Solvers

Analytic Solver V2024 Q2 also includes new (minor) releases of the **Gurobi Solver** and **OptQuest Solver** (again at no extra cost), that improve performance on a range of models using these plug-in Solver Engines. The Gurobi Solver (their version 11.0.2) has typical speed improvements of 8% to 18% compared to pre-V11 versions. The OptQuest Solver (their version 9.1.2.9) includes an improved "diversity search" algorithm, among other enhancements.

PSI Interpreter Improvements

In Analytic Solver V2024 Q2, we expect the biggest impact on performance for most users will come from some **deep** improvements in the **PSI Interpreter**, our "Polymorphic Spreadsheet Interpreter". This is the part of Analytic Solver (and RASON) that takes your model as expressed in Excel or RASON formulas, analyzes the model, and converts it into a form **usable by** the Solver Engines -- whether it's linear or quadratic (for the LSLP or Gurobi Solvers), nonlinear (for the LSGRG, LSSQP and KNITRO Solvers), or "non-smooth / arbitrary" (for the Evolutionary and OptQuest Solvers). The PSI Interpreter is hard at work when you see "Setting Up Problem..." in the Task Pane, but also **during the solution** of nonlinear and non-smooth models.

In Analytic Solver V2024 Q2, the PSI Interpreter has been re-engineered to use "sparse methods in [automatic differentiation](#)" (take a look at the Wikipedia article we've linked, if you're curious about the methods from algebra and calculus). The practical impact for users with large models is **faster** end-to-end solution times for linear models, less time spent in "**Setting Up Problem**", and big savings in **memory use** (which translates to faster solutions, or solutions for models that exhausted memory in the past).

All you really need to know is "much better performance", but if you're curious, the rest of this post will seek to explain just some of what the PSI Interpreter does for you.

More About the PSI Interpreter: The Jacobian

You might be surprised to learn that the Solver Engines, or optimization algorithms, don't work with your Excel formulas at all. They require as input *tables of numbers* (constant or changing) that describe your model, at a rather low level. The PSI Interpreter reads, parses, and interprets your formulas and produces those tables of numbers.

The most fundamental such table is called the **Jacobian matrix** -- see the [Wikipedia article](#) (rather technical) or "[A Gentle Introduction to the Jacobian](#)" which comes from machine learning, where the same matrix arises. In optimization models, this matrix has a row for the objective and each constraint, and a column for each decision variable -- so its size grows with variables *times* constraints. Each matrix element is the **partial derivative** (or *rate of change*) of one constraint (or the objective) with respect to one decision variable.

In linear models, the partial derivatives are all **constant** numbers and are referred to as "LP coefficients" -- but these numbers may not appear explicitly anywhere in your Excel model. They are computed by the PSI Interpreter. In nonlinear and non-smooth models, the partial derivatives are **not constant** -- their values depend on the current values of (potentially many) decision variables. So the PSI Interpreter must **re-compute** them, and re-supply them to the Solver Engine, each time the Solver Engine tests a new set of values for the decision variables -- and this can happen thousands or even millions of times while a model is being solved. So the PSI Interpreter has a lot of work to do.

Let's say your model has 32,000 variables and 32,000 constraints (the upper limits of our *Standard* Large-Scale LP/QP Solver -- the *Extended* version removes the limits). That's not very large by our customers' standards -- but it means the Jacobian matrix has $32,000 \times 32,000 = 1,024,000,000$ (just over **one billion**) elements. It takes a *lot* of memory, and a *lot* of computing to calculate a billion different values from your Excel model! The saving grace is that in *most* models, *most* of these matrix elements will be **zero**: A typical constraint will depend on a **small subset** of the decision variables, so its rate of change is zero with respect to all the other decision variables. The challenge for the PSI Interpreter is to **figure out which** elements will be zero (the "sparsity pattern"), **without** consuming a lot of memory and time just doing this. In V024 Q2, the PSI Interpreter is far better at this!

Beyond the Jacobian: The Hessian

When your model goes "beyond linear", even with just a quadratic objective (as in portfolio optimization), the Solver Engines will typically require a further large table of numbers, called the **Hessian matrix** (again see the [Wikipedia article](#) if you're interested). The Hessian has as many rows and as many columns as the number of **decision variables** in your model; it is typically computed for the objective function (at least). Each matrix element is the **second partial derivative** of one function (usually the objective) with respect to a **pair** of decision variables. To compute this, the PSI Interpreter must do even more work, and it must deal with the same issues of memory and computing time, for perhaps a billion values. Again *some* of these elements will be **zero**, since the function value may not depend on *all possible pairs* of variable values -- the challenge is figuring out **which ones!** And again, in V2024 Q2 we've made major progress on this. We've seen some large customer models that formerly would always exhaust memory after running for hours, suddenly solve in a minute or less!

And There's (Much) More...

There's much more to the PSI Interpreter, that we don't have time or space to fully describe here. Whenever Microsoft adds **new built-in functions** or formula features (such as "spilled arrays") to Excel, we've been hard at work supporting those new features in the PSI Interpreter. Whenever you use menu options such as "Analyze Original Model" or "Analyze Transformed Model", notice your **Model Type** (LP for linear, QP for quadratic, etc.) or the counts of model **Dependencies** in the Task Pane, or create a Linearity Report, that's the PSI Interpreter at work. And the PSI Interpreter is also the key to **super-fast Monte Carlo simulation** in Analytic Solver, usually ten times faster (or more) than other Excel add-ins for simulation. The PSI Interpreter also parses and interprets **DMN** (Decision Model and Notation) for business rules and decision tables.

You might be surprised to learn how many of our competitors offering Excel add-ins (especially firms that have been acquired by other companies) have been "resting on their laurels" for a decade or more. At Frontline Solvers, *we never rest!* We spend millions each year on real R&D, to bring you the very best in modeling and analytic software tools.

What's New in Analytic Solver V2024 Q3

In **Analytic Solver V2024 Q3**, we've gone beyond our conversational "AI Agent", introduced one year earlier in Analytic Solver V2023 Q3: Our new "AI Assistant" doesn't just consult our User Guides and publicly available info – it seeks to "understand" the *model you are working on*, so it can assist you in improving your model, or fixing problems with it!

But we don't just focus on the latest AI methods – we also improve the usability and performance of the optimization and simulation solvers you use every day.

New AI Assistant

The most visible change in V2024 Q3 is on the Analytic Solver Ribbon, which now sports two "AI icons": On the far right next to the **Help** icon is **AI Ask a Question**, similar to the capability we introduced in Analytic Solver V2023 Q3 (but enhanced): Ask it a question, and it will answer conversationally, with links to our example models, User Guides and Reference Guides. On the far left next to the **Model** icon is a new **AI Assist** icon: When you click this icon, the AI Assistant will examine and use "as context" your currently-open model in Excel (this also works in our cloud platform RASON): You can then ask a question and engage in a dialog with the Assistant about *your model*. For example, you can ask "is my model missing an essential element?" or "my model is building products from parts in inventory, how can I add another part?"

Optimization Improvements

In this release, we've included the latest versions of the **Gurobi Solver** (their 11.0.3), the **XPRESS Solver** (their 43.1.3.0), and the **MOSEK Solver** (their 10.2.0.3). For all the Solver Engines, we've expanded "Verbose" logging (enabled via the Task Pane Platform tab **General – Log Level** option) to include all the solution progress information that each Solver can provide. And in the case of the Gurobi Solver working on a non-integer, non-convex QCP

(Quadratically Constrained Problem), Verbose logging now provides “branch level” solution progress info that just wasn’t available before.

Simulation Improvements

In Monte Carlo simulation, we’ve added new “alternative parameter” distribution functions named PsiFatigueLifeAlt() and PsiFrechetAlt(). The first of these, the “fatigue life” or Birnbaum–Saunders distribution, is extensively used in reliability problems to model failure times. The second, the Frechet or inverse Weibull distribution, is used to model extreme events, for example in weather forecasting and hydroelectric power application.

What’s New in Analytic Solver V2025 Q1

We’ve just released **Analytic Solver V2025 Q1**, with new versions of Analytic Solver Desktop and Cloud (for Excel users), Solver SDK (for developers), and our cloud platform RASON. This release features both **usability** improvements and a powerful new way to solve a **wide range of nonlinear** optimization models using the **Gurobi Solver Engine** -- traditionally known for solving linear (and quadratic) mixed-integer models. This new approach yields **globally optimal** solutions **every time** -- not possible with other nonlinear Solvers -- and often **much faster** solutions.

Usability: New Flexible Dialogs

Analytic Solver has always offered an easy user interface with an Office-style Ribbon, Task Pane, and dialogs that you can move around and resize. But the **content** of many dialogs used fixed-size fonts and graphics. To help with both **accessibility**, and better use of **larger, high-resolution screens** when you’re “consuming” this content, in Analytic Solver V2025 Q1 we’ve revised the Task Pane and scores of other dialogs with “**zoomable**” content. You can now press the CTRL key, then use the mouse wheel to enlarge or shrink the text and graphics in almost any dialog or pane of Analytic Solver. At the bottom of this page is an example, where we’ve enlarged the text in the top and bottom parts of the Task Pane, but left the middle part (Model Diagnosis) unchanged – so you can see the difference.

Performance: New Nonlinear Solver Power

Analytic Solver V2025 Q1 features a *very special* level of support for the just-released (November 19) **Gurobi Solver 12.0**. This popular Solver again improves performance on linear programming (about +4%) and mixed-integer programming (about 13%) models, and even better performance (about +28%) on non-convex mixed-integer quadratic models. But most interesting for our Excel users is Gurobi’s improved support for (often non-convex) **mixed-integer nonlinear** optimization problems. Such problems arise in specialty chemicals, hydroelectric power, and some aeronautics and space applications.

Unlike other nonlinear Solvers (such as our LSGRG, LSSQP and KNITRO Solvers), Gurobi takes a different approach, that can yield **globally optimally** solutions **IF** you can formulate your model with (only) certain algebraic expressions -- defined via special Python API calls -- that use just over a [dozen mathematical operators](#) (plus, minus, times, divide, square, square root, exponential and logarithm, etc.). If you’re really good at Python coding and

really good with algebraic manipulation, you can code this yourself. But in Analytic Solver V2025 Q1, **you don't have to do any extra work** to define your nonlinear model for the Gurobi Solver.

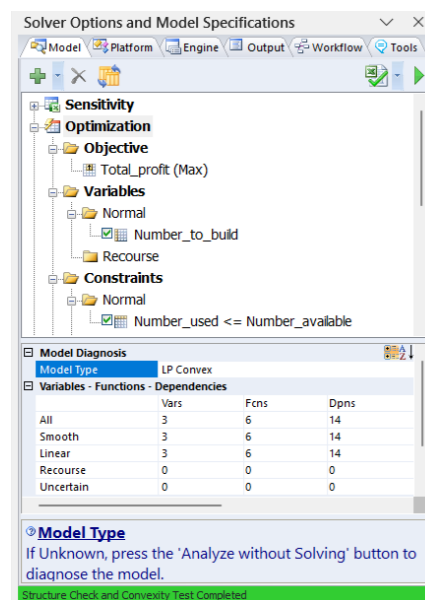
Standard Excel Formulas, New Solving Speed

You can write your model in Excel, using any of Excel's mathematical operators and any of Excel's 470+ built-in functions. Analytic Solver's **PSI Interpreter** will **analyze** your formulas (including formulas that depend on other formulas, to any level), **diagnose** your model as linear, quadratic, smooth nonlinear or non-smooth, and -- In V2025 Q1 -- tell you if your model is "nonlinear and **solvable via the new Gurobi Solver**". Analytic Solver can generate the internal "expression trees" that the Gurobi Solver requires.

Analytic Solver will even expand various Excel functions, based on their mathematical definitions, into chains of "Gurobi-supported primitive operators". For example, you can just go ahead and use the STDEV Excel function -- Analytic Solver will recognize this as "Gurobi-eligible" and create (internally) an expression tree using square, plus and square root operators, based on the definition of "standard deviation". The bottom line: With Analytic Solver, you can use the Gurobi 12.0 Solver to solve a much **wider range** of nonlinear models, much **more quickly and easily**, than you could any other way.

For "eligible" nonlinear optimization models, the Gurobi 12.0 Solver is often **faster** than other nonlinear Solvers, especially when **integer variables** are also used in the model, and when you require a **globally optimal** solution. This isn't always true however -- occasionally Gurobi is (much) slower -- so it's good that you have alternative nonlinear Solvers available to you, even in the basic Analytic Solver Optimization product. We do hope and expect to see even better performance in future versions of the Gurobi Solver.

Analytic Solver V2025 Q1 is a **free upgrade** for anyone with a paid Analytic Solver Optimization (or higher) license. If you've been a customer for years, you're probably accustomed to such dramatic enhancements -- but we're always working hard to bring you **more value, sooner** than anyone else in our "world" of advanced analytics.



Analytic Solver Product Line

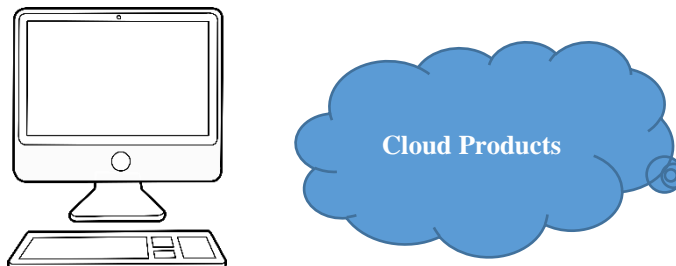
This Guide shows you how to create and evaluate forecasting and data science models using Analytic Solver.

In Frontline Solvers V2025 Q1, every license starts with **Analytic Solver Upgrade**, which allows you to **use every feature** described in this User Guide and the companion Reference Guide, for learning purposes with **small models**. Upgrade versions enable you to ‘scale up’ and solve commercial-size models for optimization, simulation or data science, paying for only what you need – but you **keep** access to all the features of Analytic Solver at the basic problem limits.

Analytic Solver combines and integrates the features of Frontline’s products for conventional optimization (formerly called **Premium Solver Pro** and **Premium Solver Platform**), Monte Carlo simulation and stochastic optimization (formerly **Risk Solver Pro** and **Risk Solver Platform**), and forecasting and data science (formerly **XLMiner Pro** and **XLMiner Platform**), in a common user interface that’s available both in Excel (desktop) and in your browser (cloud).

Analytic Solver’s **optimization** features are **fully compatible upgrades for the Solver** bundled within Microsoft Excel, which was developed by Frontline Systems for Microsoft. Your Excel Solver models and macros will work without changes; you can use either the classical **Solver Parameters dialog**, or a newer Task Pane user interface to define optimization models.

Desktop and Cloud versions



Analytic Solver V2025 Q1 includes the use of Analytic Solver Cloud and Data Science Cloud – accessible through Excel for the Web or desktop Excel. Use Analytic Solver Cloud to solve optimization and simulation models and use the Data Science Cloud app to perform data science or forecasting. Both the Analytic Solver and Data Science Cloud apps support existing models created in previous versions of Analytic Solver. Your license for Analytic Solver will allow you to use Analytic Solver Desktop in desktop Excel or Analytic Solver Cloud in either desktop Excel¹ or Excel for the Web. A license for a desktop product will grant you a license for the corresponding product in a cloud app. For example, a license for Analytic Solver Optimization in Analytic Solver Desktop will also grant you a license for Analytic Solver Optimization in the Analytic Solver Cloud app. The overwhelming majority of features in Analytic Solver Desktop are also included in Analytic Solver Cloud and Data Science

¹ Must be using Microsoft Excel V2020 or later.

Cloud. However, there could be slight differences in the way users execute a function in the cloud apps.

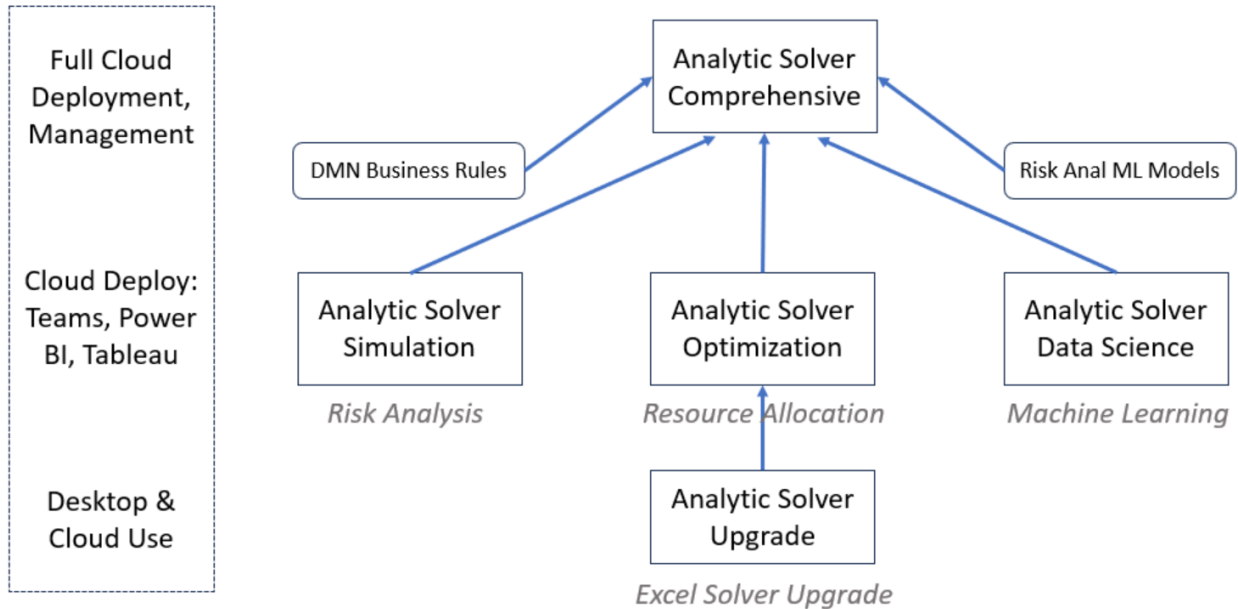
Analytic Solver for desktop or cloud may be purchased in several different ways starting with the most basic version, Analytic Solver Upgrade, up to our most complete version, Analytic Solver Comprehensive. Continue reading to see which product will best meet your needs.

Analytic Solver Academy

Analytic Solver Academy allows access to Frontline Systems' Solver Academy courses: Optimization 101, Simulation 101, Data Mining 101 and Data Mining 201. These four courses will quickly bring users up to speed on the complete functionality of Analytic Solver. If no other license exists, this product offers a basic course license with size limits for Monte Carlo simulation, data science and text mining sufficient to run all installed example models and complete all course work.

Product Licenses

Frontline Systems offers six annual licenses for our Desktop Excel and Excel Online product suite: Analytic Solver Upgrade, Analytic Solver Optimization, Analytic Solver Simulation, Analytic Solver Data Mining, and Analytic Solver Comprehensive. The graphic below illustrates how each of the subsets roll into Analytic Solver Comprehensive, our flagship product.



Analytic Solver Upgrade

Frontline Systems is Frontline's basic upgrade for the Excel Solver – enabling you to solve linear up to 10 times larger (up to 2,000 variables), quadratic (QPs) and quadratically constrained problems (QCPs) and nonlinear models 2.5 times larger (up to 500 variables). It includes faster versions of the LP/Quadratic, GRG Nonlinear, and Evolutionary Solvers, but it doesn't support plug-in large-scale Solver Engines, decision tables, box functions, model deployment or the

AI Agent. This product includes basic level support for simulation and data science models.

Analytic Solver Optimization, Simulation, Data Science and Comprehensive

Users of the next three products are granted a license to share the results of their optimization or simulation model to Power BI, Tableau or Microsoft Teams via the Deploy Model button on the Analytic Solver tab. See the *Deploying Your Model* chapter within the Analytic Solver User Guide for more details.

Analytic Solver Optimization

Analytic Solver Optimization (formerly **Premium Solver Platform**) is Frontline's most powerful product for *conventional* optimization. It includes the PSI Interpreter, five built-in Solvers (LP/Quadratic, SOCP Barrier, GRG Nonlinear, Interval Global, and Evolutionary), solves linear, quadratic and quadratically constrained models up to 8,000 variables and nonlinear models up to 1,000 variables, and it supports plug-in large-scale Solver Engines to handle much larger models. When used with Analytic Solver Simulation, you can also solve models *with uncertainty* using simulation optimization, stochastic linear programming, and robust optimization. This product, alone, includes all the features of Analytic Solver Upgrade with limited size constraints.

Analytic Solver Simulation

Analytic Solver Simulation (*expanded* from **Risk Solver Pro**) is Frontline's full-function product for Monte Carlo simulation and simulation optimization. It includes decision tree capabilities and the PSI Interpreter – which gives you the fastest Monte Carlo simulations available in any Excel-based products, unique interactive simulation capabilities, multiple parameterized simulations, and simulation optimization using the Evolutionary Solver. When coupled with Analytic Solver Optimization, you can also solve models with uncertainty using stochastic linear programming and robust optimization. This product, alone, includes all the features of Analytic Solver Upgrade with limited size constraints.

Automated Risk Analysis of Machine Learning Models

Analytic Solver Data Science and Analytic Solver Comprehensive offer capability for **automated risk analysis** of **machine learning** models along with **synthetic data generation**, to augment the data you already have. The beauty of these new features is that **you don't have to do any work** to obtain a risk analysis of your model's performance, beyond a few mouse clicks on the **Simulation** tab to enable the analysis, which is entirely **automated**. You don't even have to be familiar with the features of Analytic Solver Simulation to use this capability for machine learning!

Analytic Solver Data Science

Analytic Solver Data Science (formerly XLMiner Platform and more recently Analytic Solver Data Mining) is Frontline's most powerful product for data science, text mining, forecasting and predictive analytics. It includes data access and sampling, data exploration and visualization, text mining, data transformation, and feature selection capabilities; time series forecasting with

ARIMA and exponential smoothing; and a wide range of data science methods for classification, prediction and affinity analysis, from multiple regression to neural networks. This product includes all the features of Analytic Solver Upgrade.

Analytic Solver Comprehensive

Analytic Solver Comprehensive (formerly **Analytic Solver Platform**) combines the optimization capabilities of Analytic Solver Optimization, the simulation capabilities of Analytic Solver Simulation, and the data science capabilities of Analytic Solver Data Science (formerly XLMiner and more recently Analytic Solver Data Mining). It includes the PSI Interpreter, five built-in Solvers (LP/Quadratic, SOCP Barrier, GRG Nonlinear, Interval Global, and Evolutionary) and it accepts a full range of plug-in large-scale Solver Engines. It supports optimization, Monte Carlo simulation, simulation optimization, stochastic programming and robust optimization, and large-scale data science and forecasting capabilities.

See the chart below to compare the data handling limits between Analytic Solver Data Science and the remaining subsets: Analytic Solver Optimization/Simulation. (Analytic Solver is available for classroom use or by purchasing a textbook that contains the software. For a complete list of textbooks that include Analytic Solver, see our Website at www.solver.com.) "Unlimited" indicates that no limit is imposed, however, other limits may apply based on your computer resources and Excel version.

bn	Data Science/Comprehensive	Subsets
Risk Analysis and Synthetic Data Generation	Unlimited	Not supported
Partitioning		
# of Records (original data)	Unlimited	65,000
# of Records (training partition)	Unlimited	10,000
# of Variables (output)	Unlimited	50
Sampling		
# of Records (original data)	Unlimited	65,000
# of Variables (output)	Unlimited	50
# of Strata (Stratified Sampling)	Unlimited	30
Database		
# of Records (table)	Unlimited	1,000,000
# of Records (output)	Unlimited	65,000
# of Variables (table)	Unlimited	Unlimited
# of Variables (output)	Unlimited	50
# of Strata (Stratified Sampling)	Unlimited	30
File System		
# of Files	Unlimited	100
Text Mining		
# Documents	Unlimited	100
# Characters (per document)	Unlimited	5,000
# Terms in final vocabulary	Unlimited	50
# Text columns	Unlimited	1

Transformation		
<i>Common</i>		
# Records	Unlimited	10,000
# Variables	Unlimited	50
Missing Data Handling		
# of Records	Unlimited	65,000
# of Variables	Unlimited	50
Binning Continuous Data		
# of Records	Unlimited	65,000
Transforming Categorical Data		
# of Records	Unlimited	65,000
# of Variables (data range)	Unlimited	50
# of Distinct values	Unlimited	30
Time Series Analysis		
# of Records	Unlimited	1,000
Classification and Prediction		
# of Records (total)	Unlimited	65,000
# of Records (training partition)	Unlimited	10,000
# of Records (new data for scoring)	Unlimited	65,000
# of Variables (output)	Unlimited	50
# of Distinct classes (output variable)	Unlimited	30
# of Distinct values (categorical input variables)	Unlimited	15
k-Nearest Neighbors		
# of Nearest neighbors	50	10
Regression/Classification Trees		
# of Splits	Unlimited	100
# of Nodes	Unlimited	100
# of Levels	Unlimited	100
# Levels in Tree Drawing	Unlimited	7
Ensemble Methods		
# Weak learners	Unlimited	10
Feature Selection		
# of Records	Unlimited	10,000
# of Variables	Unlimited	50
# of Distinct classes (output variable)	Unlimited	30
# of Distinct values (input variables)	Unlimited	100
Association Rules		
# of Transactions	Unlimited	65,000
# of Distinct items	5,000	100
Clustering		
<i>K-Means</i>		

# of Records	Unlimited	10,000
# of Variables	Unlimited	50
# of Clusters	Unlimited	10
# of Iterations	Unlimited	50
<i>Hierarchical</i>		
# of Records	Unlimited	10,000
# of Variables	Unlimited	50
# of Clusters in a Dendrogram	Unlimited	10
Size of distance matrix	Unlimited	1,000 x 1,000
Charts		
# of Records	Unlimited	65,000
# of Variables (original data)	Unlimited	100
General		
Model pane	Included	Included
Big Data sampling/summarization	Included	Included
Model storage and scoring	Included	Included

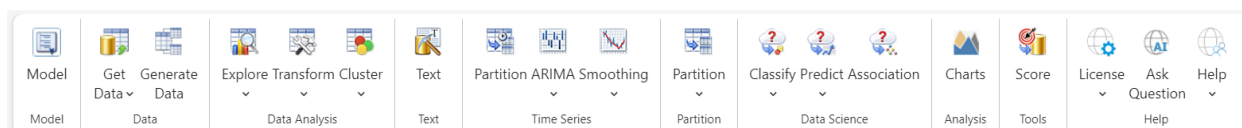
Data Science Ribbon Overview

Analytic Solver Data Science software offers over 30 different methods for analyzing a dataset in order to forecast future events. The Data Science ribbon is broken up into five different segments as shown in the screenshot below.

Desktop Analytic Solver Data Science



Data Science Cloud



- Click the **Model** button to display the Solver Task Pane. This new feature (added in V2016) allows you to quickly navigate through datasets and worksheets containing Analytic Solver Data Science results.
- Click the **Get Data** button to draw a random sample of data, or summarize data from a (i) an Excel worksheet, (ii) the PowerPivot “spreadsheet data model” which can hold 10 to 100 million rows of data in Excel, (iii) an external SQL database such as Oracle, DB2 or SQL Server, or (iv) a dataset with up to billions of rows, stored across many hard disks in an external **Big Data** compute cluster running Apache Spark (<https://spark.apache.org/>).
- You can use the **Data Analysis** group of buttons to explore your data, both visually and through methods like cluster analysis, transform your data with methods like Principal Components, Missing Value imputation, Binning

continuous data, and Transforming categorical data, or use the Text Mining feature to extract information from text documents.

- Use the **Time Series** group of buttons for time series forecasting, using both Exponential Smoothing (including Holt-Winters) and ARIMA (Auto-Regressive Integrated Moving Average) models, the two most popular time series forecasting methods from classical statistics. These methods forecast a single data series forward in time.
- The **Data Science** group of buttons give you access to a broad range of methods for prediction, classification and affinity analysis, from both classical statistics and data science. These methods use multiple input variables to predict an outcome variable or classify the outcome into one of several categories. Introduced in V2015, Analytic Solver Data Science and now the Data Science Cloud app, offer Ensemble Methods for use with Classification Trees, Regression Trees, and Neural Networks.
- Use the **Predict** button to build prediction models using Multiple Linear Regression (with variable subset selection and diagnostics), k-Nearest Neighbors, Regression Trees, and Neural Networks. Use Ensemble Methods with Regression Trees and Neural Networks to create more accurate prediction models.
- Use the **Classify** button to build classification models with Discriminant Analysis, Logistic Regression, k-Nearest Neighbors, Classification Trees, Naïve Bayes, and Neural Networks. Use Ensemble Methods with Classification Trees and Neural Networks to create more accurate classification models.
- Use the **Associate** button to perform affinity analysis (“what goes with what” or market basket analysis) using Association Rules.
- Use the **Score** button to score new data using a fitted forecasting, classification or prediction model.
- If forecasting and data science are new for you, don’t worry – you can learn a lot about them by consulting our **AI Agent**, Frontline’s artificial intelligence technical support assistant. AI Agent is designed to provide assistance and support for users of Frontline Solvers’ Analytic Solver and Analytic Solver Data Science software. The AI Agent is knowledgeable about the functionality and features of the software, as well as the concepts and processes involved in optimization, simulation and data science/forecasting. Just enter a topic or question such as “What classification algorithms are supported in Analytic Solver Data Science?” and click Submit Query to get started.
- Use the **License** button to manage your account and licenses.
- Use the **Help** button to open example models, open the Help Center, where you can find pre-recorded webinars or access our Knowledge Base or explore our User Guides.

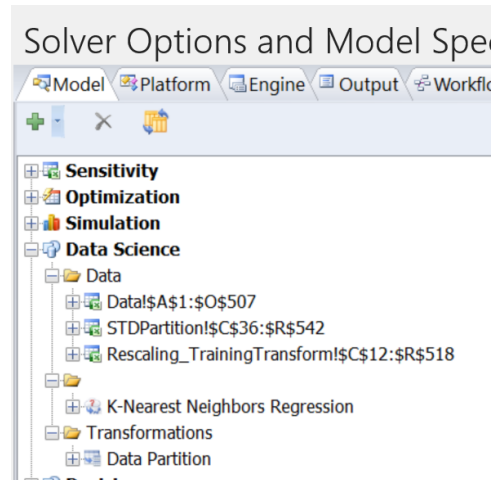
If you’d like to learn more and get started as a ‘data scientist,’ consult the excellent book *Data Mining for Business Intelligence*, which was written by the original Data Science (formally known as XLMiner and most recently Analytic Solver Data Mining) designers and early academic users. You’ll be able to run all the Data Science examples and exercises in Analytic Solver.

Analytic Solver Data Science, along with the Data Science Cloud app, can be purchased as a stand-alone product. A stand-alone license for Analytic Solver Data Science includes all of the data analysis, time series data capabilities,

classification and prediction features available in Analytic Solver Comprehensive but does not support optimization or simulation. See the Analytic Solver Data Science User Guide Data Specifications for each product.

Model

Click the Model button in Analytic Solver Desktop to display the Solver Task Pane within Analytic Solver Data Science. From the Model tab, you can easily navigate between worksheets containing data and results. Note that in the Data Science Cloud app, only the Workflow tab is present on the task pane.



All fields contained in the dataset are listed under the name of each data containing worksheet (for example Data!\$A\$1:\$O\$507) while results obtained from Analytic Solver Data Science are listed under Reports or Transformations by type and run number.

See the next chapter, Creating Workflows for information on how to create a workflow in Analytic Solver Data Science.

Get Data

Analytic Solver Data Science includes several different methods for importing your data including Sampling from either a Worksheet or Database or Importing from a File Folder.



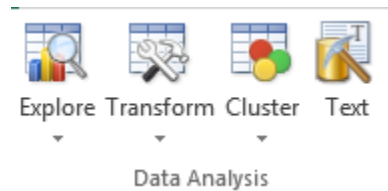
Click the **Get Data** icon to take a representative sample from a dataset included in either an Excel workbook or an Oracle, SQL Server, MS-Access, or Power Pivot database. Users can choose to sample with or without replacement using simple or stratified random sampling. Click **Get Data – File Folder** to import and or sample from a collection of text documents for use with Text Mining.

Click the **Get Data - Big Data** to sample from or summarize from a dataset with up to billions of rows, stored across many hard disks in an external compute

cluster running Apache Spark. Results may be requested immediately upon completion or at a later time using **Get Results**.

Data Analysis

Analytic Solver Data Science includes several different methods for data analysis, including Sampling from either a Worksheet or Database (not supported in the Cloud app), Charting with 8 different types of available charts, Transformation techniques which handle missing data, binning continuous data, creating dummy variables and transforming categorical data, and using Principal Components Analysis to reduce and eliminate superfluous or redundant variables; along with two different types of Clustering techniques, k-Means and Hierarchical.



Click the Explore icon to use Feature Selection to help decide which variables should be included in your classification or prediction models or use the Chart Wizard to create one or more charts of your data. The new Feature Selection tool can help give insight into which variables are the most important or relevant for inclusion in your classification or prediction model using various types of statistics and data analysis measures. Analytic Solver Data Science includes 8 different types of charts to choose from, including: bar charts, line charts, scatterplots, boxplots, histograms, parallel coordinates charts, scatterplot matrix charts or variable charts. Click this icon to edit or view previously created charts as well.

Click the Transformation icon when data manipulation is required. In most large databases or datasets, a portion of variables are bound to be missing some data. Analytic Solver Data Science includes routines for dealing with these missing values by allowing a user to either delete the full record or apply a value of her/his choice. Analytic Solver Data Science also includes a routine for binning continuous data for use with prediction and classification methods which do not support continuous data. Continuous variables can be binned using several different user specified options. Non-numeric data can be transformed using dummy variables with up to 30 distinct values. If more than 30 categories exist for a single variable, use the Reduce Categories routine to decrease the number of categories to 30. Finally, use Principal Components Analysis to remove highly correlated or superfluous variables from large databases.

Click the Cluster icon to gain access to two different types of clustering techniques: k-Means clustering and hierarchical clustering. Both methods allow insight into a database or dataset by performing a cluster analysis. This type of analysis can be used to obtain the degree of similarity (or dissimilarity) between the individual objects being clustered.

Click the Text icon to use the Text Miner tool to analyze a collection of text documents for patterns and trends. (In the Cloud app, this tool is included in the Text section of the Ribbon.) These algorithms can categorize documents, provide links between documents that were not otherwise noted and create visual maps of the documents. Analytic Solver Data Science takes an integrated approach to text mining by combining text processing and analysis in a single

package. While Analytic Solver Data Science is effective for mining “pure text” such as a set of documents, it is especially useful for “integrated text and data science” applications such as maintenance reports, evaluation forms, or any situation where a combination of structured data and free-form text data is available.

Time Series Analysis

Analytic Solver Data Science also supports the analysis and forecasting of datasets that contain observations generated sequentially such as predicting next year’s sales figures, monthly airline bookings, etc. through partitioning, autocorrelations or ARIMA models and through smoothing techniques.

A time series model is first used to obtain an understanding of the underlying forces and structure that produced the data and then secondly, to fit a model that will predict future behavior. In the first step, the analysis of the data, a model is created to uncover seasonal patterns or trends in the data, for example bathing suit sales in June. In the second step, forecasting, the model is used to predict the value of the data in the future, for example, next year's bathing suit sales. Separate modeling methods are required to create each type of model.



Typically, when using a time series dataset, the data is first partitioned into training and validation sets. Click the Partition icon within the Time Series ribbon segment to utilize the Time Series Data Partitioning routine. Analytic Solver Data Science features two techniques for exploring trends in a dataset, ACF (Autocorrelation function) and PACF (Partial autocorrelation function). These techniques help the user to explore various patterns in the data which can be used in the creation of the model. After the data is analyzed, a model can be fit to the data using the ARIMA method. All three of these methods can be found by clicking the ARIMA icon on the Data Science ribbon.

Data collected over time is likely to show some form of random variation. "Smoothing techniques" can be used to reduce or cancel the effect of these variations. These techniques, when properly applied, will “smooth” out the random variation in the time series data to reveal any underlying trends that may exist.

Click the Smoothing icon to gain access to Analytic Solver Data Science’s four different smoothing techniques: Exponential, Moving Average, Double Exponential, and Holt Winters. The first two techniques, Exponential and Moving Average, are relatively simple smoothing techniques and should not be performed on datasets involving trends or seasonality. The third technique, Double Exponential, should be used when a trend is present in the dataset, but not seasonality. The last technique, Holt Winters, is a more advanced technique and should be selected when working with datasets involving seasonality.

Data Science

The Data Science section of the Data Science ribbon contains four icons: Partition, Classify, Predict, and Associate. Note that in the Cloud app, Partition has been cordoned off into its own "Partition" section. Click the Partition icon to partition your data into training, validation, and if desired, test sets. Click the Classify icon to select one of six different classification methods. Click the Predict icon to select one of four different prediction methods. Click the Associate icon to recognize associations or correlations among variables in the dataset.

Analytic Solver Data Science supports six different methods for predicting the class of an outcome variable (classification) along with three ensemble methods which use these six methods as weak learners, and four different methods, along with three ensemble methods, for predicting the actual (prediction) of an outcome variable. Classification can be described as categorizing a set of observations into predefined classes in order to determine the class of an observation based on a set of variables. A prediction method can be described as a technique performed on a database either to predict the response variable value based on a predictor variable or to study the relationship between the response variable and the predictor variables. For example, when determining the relationship between the crime rate of a city or neighborhood and demographic factors such as population, education, male to female ratio, etc.

One very important issue when fitting a model is how well the newly created model will behave when applied to new data. To address this issue, the dataset can be divided into multiple partitions before a classification or prediction algorithm is applied: a training partition used to create the model, a validation partition to test the performance of the model and, if desired, a third test partition. Partitioning is performed randomly, to protect against a biased partition, according to proportions specified by the user or according to rules concerning the dataset type. For example, when creating a time series forecast, data is partitioned by chronological order.

The six different classification methods are:

Discriminant Analysis - Constructs a set of linear functions of the predictor variables and uses these functions to predict the class of a new observation with an unknown class. Common uses of this method include: classifying loan, credit card or insurance applicants into low or high risk categories, classifying student applications for college entrance, classifying cancer patients into clinical studies, etc.

Logistic Regression – A variant of ordinary regression which is used to predict the response variable, or the output variable, when the response variable is a dichotomous variable (a variable that takes only two values such as yes/no, success/failure, survive/die, etc.).

k-Nearest Neighbors – This classification method divides a training dataset into groups of k observations using a Euclidean Distance measure to determine similarity between “neighbors”. These classification groups are used to assign categories to each member of the validation training set.

Classification Tree – Also known as Decision Trees, this classification method is a good choice when goal is to generate easily understood and explained “rules” that can be translated in an SQL or query language.

Naive Bayes – This classification method first scans the training dataset and finds all records where the predictor values are equal. Then the most prevalent class of the group is determined and assigned to the entire collection of observations. If a new observation's predictor variable equals the predictor variable of this group, the new observation will be assigned to this class. Due to the simplicity of this method a large number of records are required to obtain accuracy.

Neural Network – Artificial neural networks are based on the operation and structure of the human brain. These networks process one record at a time and “learn” by comparing their classification of the record (which as the beginning is largely arbitrary) with the known actual classification of the record. Errors from the initial classification of the first records are fed back into the network and used to modify the networks algorithm the second time around. This continues for many, many iterations.

The four different predictive methods are:

Multiple Linear Regression – This method is performed on a dataset to predict the response variable based on a predictor variable or used to study the relationship between a response and predictor variable, for example, student test scores compared to demographic information such as income, education of parents, etc.

k-Nearest Neighbors – Like the classification method with the same name above, this prediction method divides a training dataset into groups of k observations using a Euclidean Distance measure to determine similarity between “neighbors”. These groups are used to predict the value of the response for each member of the validation set.

Regression Trees - A Regression tree may be considered a variant of a decision tree, designed to approximate real-valued functions instead of being used for classification methods. As with all regression techniques, Analytic Solver Data Science assumes the existence of a single output (response) variable and one or more input (predictor) variables. The output variable is numerical. The general regression tree building methodology allows input variables to be a mixture of continuous and categorical variables. A decision tree is generated when each decision node in the tree contains a test on some input variable's value. The terminal nodes of the tree contain the predicted output variable values.

Neural Network – Artificial neural networks are based on the operation and structure of the human brain. These networks process one record at a time and “learn” by comparing their prediction of the record (which as the beginning is largely arbitrary) with the known actual value of the response variable. Errors from the initial prediction of the first records are fed back into the network and used to modify the networks algorithm the second time around. This continues for many, many iterations.

Three ensemble methods, bagging, boosting, and random trees, are also available for both classification and prediction. Each of these methods uses a classification or prediction method as a weak learner. Each of these can be accessed on the button of either the Classify or Predict menus.

The goal of association rule mining is to recognize associations and/or correlations among large sets of data items. A typical and widely-used example of association rule mining is the Market Basket Analysis. Most 'market basket'

databases consist of a large number of transaction records where each record lists all items purchased by a customer during a trip through the check-out line. Data is easily and accurately collected through the bar-code scanners. Supermarket managers are interested in determining what foods customers purchase together, like, for instance, bread and milk, bacon and eggs, wine and cheese, etc. This information is useful in planning store layouts (placing items optimally with respect to each other), cross-selling promotions, coupon offers, etc.

Scoring

Click the Score icon to score new data in a database or worksheet with any of the Classification or Prediction algorithms. This facility matches the input variables to the database (or worksheet) fields and then performs the scoring on the database (or worksheet).

Analytic Solver Data Science also supports the scoring of Test Data. When Analytic Solver Data Science calculates prediction or classification results, internal values and coefficients are generated and used in the computations. Analytic Solver Data Science saves these values to an additional output sheet, termed Stored Model Sheet, which uses the output sheet name, `XX_Stored_N` where `XX` are the initials of the classification or prediction method and `N` is the number of generated stored sheets. This sheet is used when scoring the test data.

Note: In previous versions of XLMiner, this utility was a separate add-on application named XLMLCalc. Starting in XLMiner V12.5, this utility is included free of charge. Starting in V2014-R2, `PsiClassify()`, `PsiPredict()` and `PsiForecast` functions are available for instantaneous interactive scoring on the worksheet without the need to click the Score icon.

Using Help, Licensing and Product Subsets

Introduction

Analytic Solver Data Science (previously referred to as XLMiner™ and more recently Analytic Solver Data Mining) is a comprehensive data science software package for use on the Web or as an add-in to Excel. Data science is a discovery-driven data analysis technology used for identifying patterns and relationships in data sets. With overwhelming amounts of data now available from transaction systems and external data sources, organizations are presented with increasing opportunities to understand their data and gain insights into it. Data science is still an emerging field, and is a convergence of fields like statistics, machine learning, and artificial intelligence.

Often, there may be more than one approach to a problem. Analytic Solver Data Science is a tool belt to help you get started quickly by offering a variety of methods to analyze your data. It has extensive coverage of statistical and machine learning techniques for classification, prediction, affinity analysis and data exploration and reduction.

This chapter describes the ways Analytic Solver V2018 handles the overall operation, including registration, licensing, use of product subsets, and use of the Startup Screen, online Help and examples.

Working with Licenses in V2025 Q1

A **license** is a grant of rights, from Frontline Systems to you, to use our software in specified ways. Information about a license – for example, its temporary vs. permanent status and its expiration date – is encoded in a **license code**. The same binary files are used for all Analytic Solver products. The product features you see depend on the license code you have.

Frontline License Manager:

In Analytic Solver V2025 Q1, the Frontline License Manager has replaced the Reprise License Manager – its basic purpose is license control (allowing the software to run, or not). But unlike Reprise, the new Frontline License Manager ties a license to a human user ID / email, not to a hardware ID or “lock code”. Starting with a user ID stored locally, it will ask for license rights for a user via a web (REST API) request to Frontline's License Manager, and store this license information locally. When our software is first installed, it will have an embedded user ID *and* license code for a trial or evaluation license, so the software can run for the trial period even without Internet access.

Upgrading from an Early Version of Desktop ASP

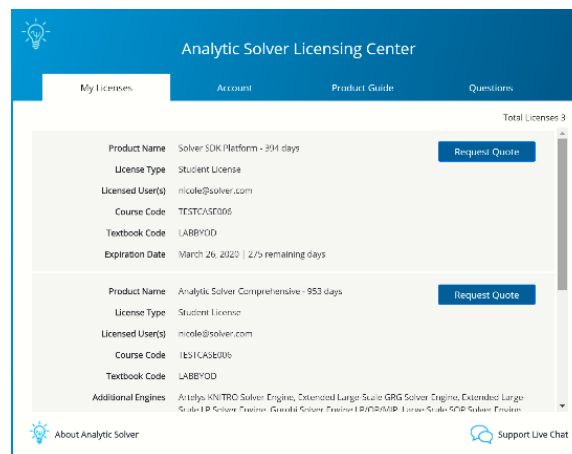
Analytic Solver Desktop V2025 Q1 can store license codes locally but typically your license will be stored on Frontline's license server. If you're upgrading from V2017/2018, a *new* V2025 Q1 license code will **not** be required. Simply

download and install V2025 Q1 for Desktop Excel and your existing V2017/V2018 license will activate the newer version. Old license codes for V2016 and earlier have no negative effect in V2025 Q1. If they exist in the obsolete Solver.lic file (located at C:\ProgramData\Frontline Systems), they will be ignored. A license code will be issued at the time of purchase.

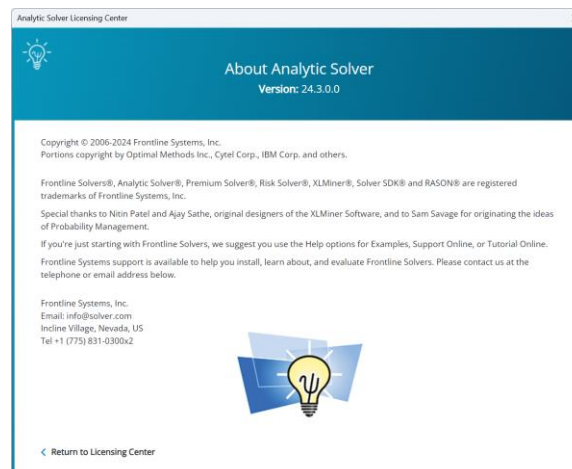
Managing Your Licenses

Click the License button to open the License Manager where you can manage your current licenses and accounts, open our Product Selection Wizard, connect to Live Chat or peruse through a list of FAQs.

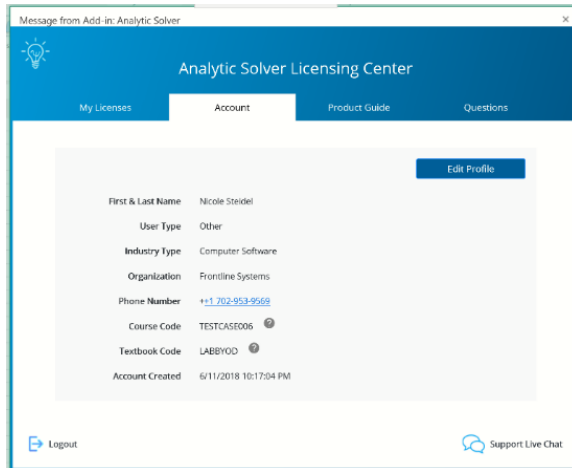
The "MyLicenses" tab displays your current license and license type, along with the expiration date. You can request a quote to renew your current license or, if your license has expired or is within 30 days of expiring, you can purchase a new license through our online store.



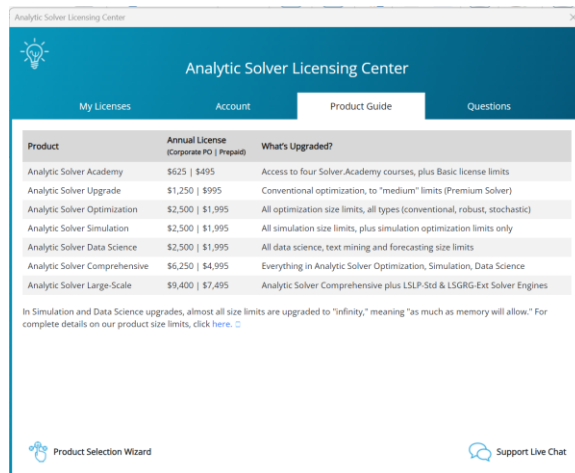
Click *About Analytic Solver* to open the following dialog containing information on this release.



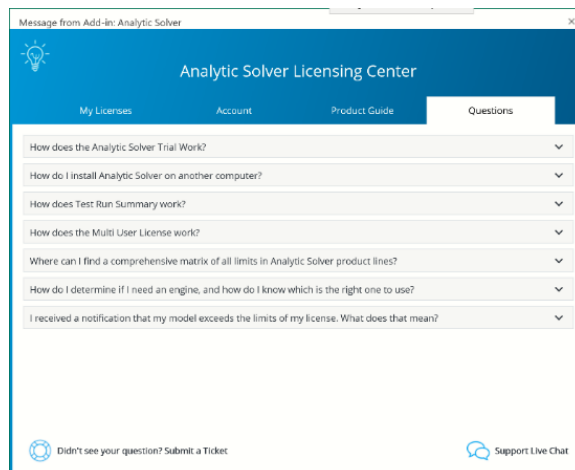
Click the Account tab to view your account on www.solver.com. Click Edit Profile to edit the information. Click Live Chat to open a Live Chat window or Log Out to log out of the product.



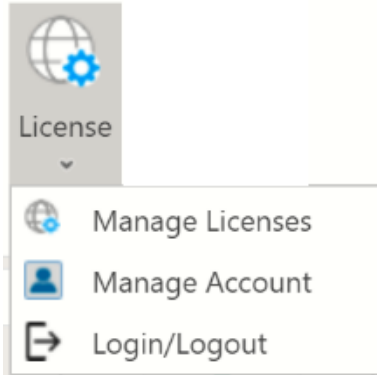
Click the Product Guide tab to view a list of products and pricing information. Click Product Selection Wizard to open the Product Selection Wizard. See the next section for information on this feature.



Click the Questions tab to review a list of FAQs, submit a support ticket or start a live chat.



Use the License menu to gain shortcuts to your account and to login or logout of Analytic Solver.



Product Selection Wizard

Select **Product Selection Wizard** from the Product Guide tab in the Licensing Center to open a series of dialogs that will help you determine which product will best meet your needs based on your recent pattern of use.

ProductSelectionWizard

Welcome to the **Product Selection Wizard!** Since you can use – and pay for – only what you need, this Wizard will help you choose from the available license options.

Analytic Solver's features cover three main problem solving areas – what do you want to do in each area?

Analytic Area	I want to gain modeling skills, or build a proposal/prototype	I have a current project to build a significant model of this type
Optimization	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Simulation/Risk Analysis	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Data Science	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

With any paid Analytic Solver license, you can always use **all** optimization, simulation, and data science features to build small models! But licenses have different **size limits** on models and data.

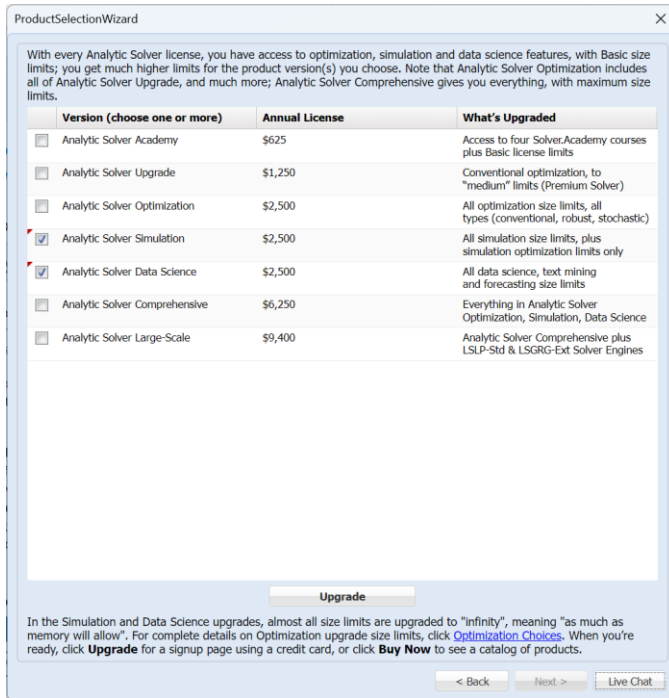
The Optimization upgrade you need depends on your model **type** (linear, nonlinear, integer), **size** and **complexity**. There are three basic levels of Optimization upgrades:

Optimization License Upgrade	Annual License
<input type="checkbox"/> Analytic Solver Upgrade (formerly Premium Solver Pro)	Contact us for details
<input type="checkbox"/> Analytic Solver Optimization (formerly Premium Solver Platform)	Contact us for details
<input type="checkbox"/> Analytic Solver Optimization + plug-in Solver Engine (Analytic Solver Large-Scale offers special discount)	Contact us for details

For more details on specific size limits enabled by these upgrades, click [Optimization Choices](#).

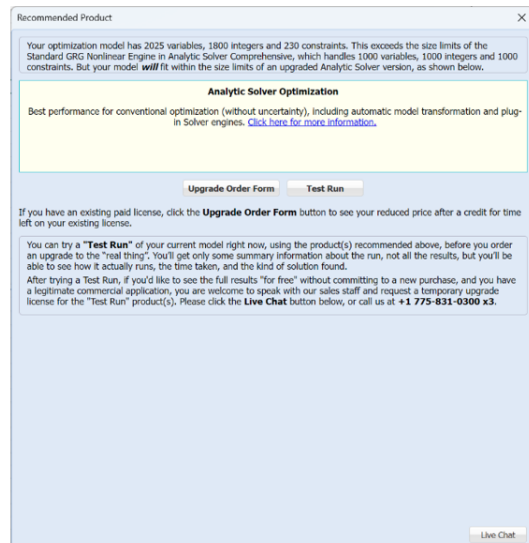
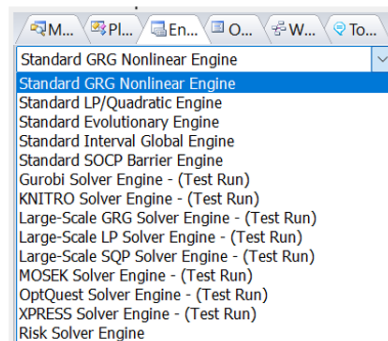
< Back Next > Live Chat

Select the Product that you'd like to purchase and then click **Next**. Click the *Optimization Choices* link to learn more about Analytic Solver products that can solve optimization models and to find more information on speed, memory, and the use of plug-in Solver Engines.

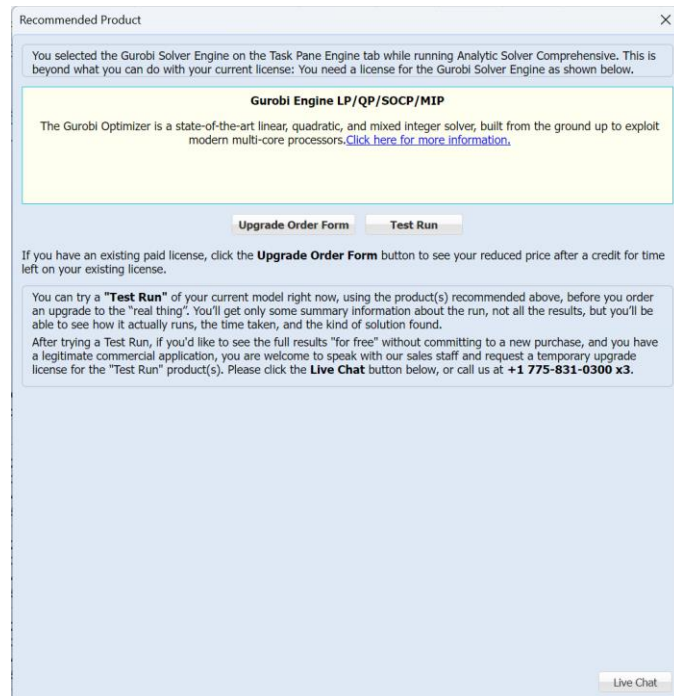


On this screen, the Product Selection Wizard will recommend a product or products based on your answers on the previous screens. Click Upgrade to purchase the recommended product. Click the *Optimization Choices* link to learn more about Analytic Solver products that can solve optimization models. If at any time you'd like to chat with a member of our Technical Support staff, click **Live Chat**. Or if you'd like to amend your answers on a previous dialog, click **Back**.

When you run a simulation or optimization model that contains too many decision variables/uncertain variables or constraints/uncertain functions for the selected engine, the Product Wizard will automatically appear and recommend a product that *can* solve your model.



When you click “Test Run”, the Product Wizard will immediately run the optimization or simulation model using the recommended product. (Only summary information will be available.) At this point, you can purchase the recommended product(s), or close the dialog.



This same behavior will also occur when solving smaller models, if you select a specific external engine, from the Engine drop down menu on the Engine tab of the Solver Task Pane, for which you do not have a license. The Product Wizard will recommend the selected engine, and allow you to solve your model using this engine. Once Solver has finished solving, you will have the option to purchase the product.

Getting Help

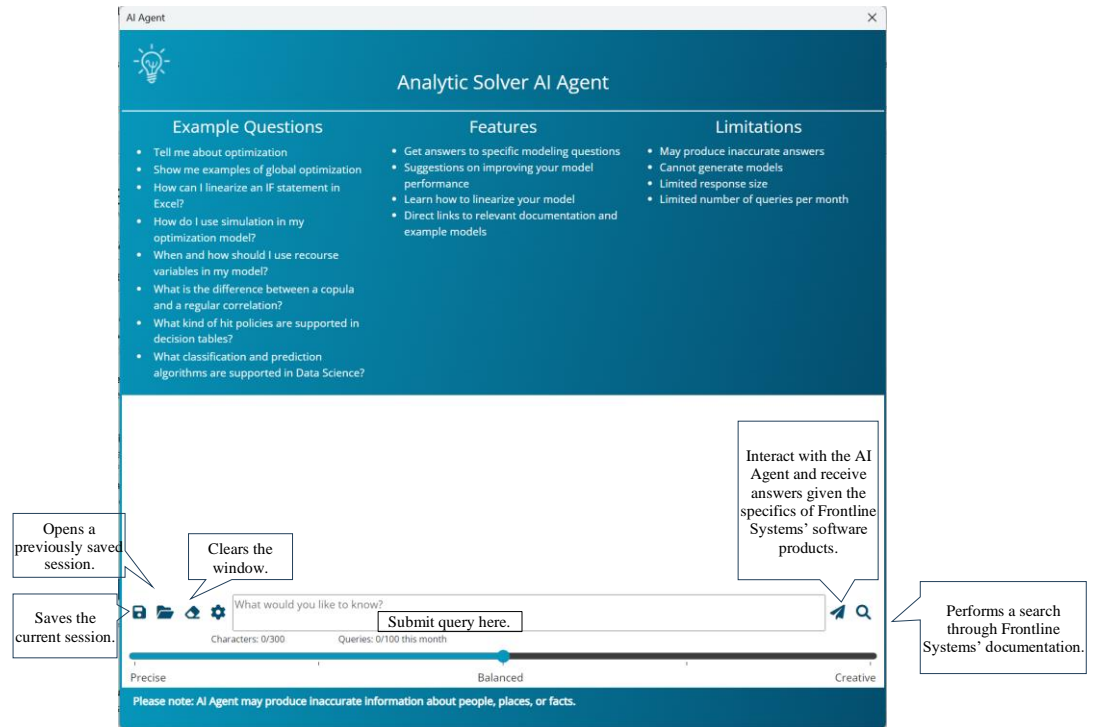
The AI Help features in Analytic Solver are designed to help users of Frontline Systems' software, including Analytic Solver and Analytic Solver Data Science, by providing technical support and guidance.

- Click AI Assist (on the far left of the Analytic Solver ribbon) while your model is open to allow AI Assist to analyze your model and check the formulation for correctness.
- Click Ask Question (on the far right of the Analytic Solver and Data Science ribbons) to open the AI Agent. AI Agent is knowledgeable about the functionalities and features of the software, as well as the underlying concepts and processes involved in optimization, simulation, forecasting, and data science.

Ask Question

If you are new to Analytic Solver and/or optimization, simulation, forecasting and data science, don't worry – Frontline's AI technical support assistant, AI

Agent, is here to help. AI Agent is designed to provide assistance and support for users of Frontline Solvers' Analytic Solver and Analytic Solver Data Science software. The AI Agent is knowledgeable about the functionality and features of the software, as well as the concepts and processes involved in optimization, simulation and data science/forecasting. Just enter a topic or question such as “What classification algorithms are supported in Analytic Solver Data Science?” and click Submit Query to get started.



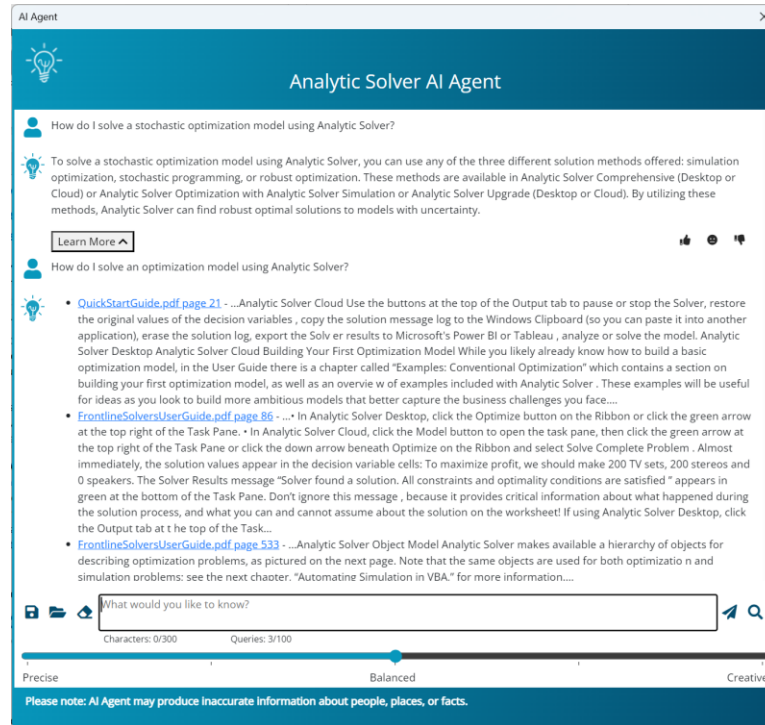
Toggle between Precise, Balanced, and Creative to determine the type of answers returned. If Precise is selected, AI Agent will attempt to be as exact and deterministic as possible while Creative will usually result in more original, uncertain and non-repetitive answers. Use Balanced (the default) for the best of both worlds.

Click the **Save** file icon to save the current session. The file, transcript.json, will be downloaded. Click the **Open** file icon to open and restore a previously saved session. Click the **Erase** icon, to clear the current search window and start a new topic.

After typing a query,

- Click **Submit Query** to interact with the AI Agent and receive answers given the specifics of Frontline Systems' software products. On each new question, the AI Agent will take into account the chat history that is present in the current session. Click *Learn More* to view all extracts of source documents that are most relevant to the query. Each extract contains a link to the appropriate page in one of the official documents.
- Click **Search Documentation** (for faster search results) to perform a semantic/relevance search through all Analytic Solver and RASON documentation. The resulting response will be similar to the “Learn More” section as described above.

The screenshot below illustrates the difference between the two different types of queries. The first query used AI technology by using Submit Query. The 2nd query was performed using Search Documentation.



Be sure to rate answers to improve future performance!

AI Assist

Frontline Systems' new AI assist is a feature designed to enhance model formulation for optimization, simulation, and stochastic models. AI Assist helps users by:

1. **Analyzing Model Formulation:** AI Assist incorporates artificial intelligence to analyze your existing model and giving you tips on how to expand and improve your current model formulation.
2. **Providing Recommendations:** The AI can offer suggestions for improving model performance or adjusting parameters to achieve better results.
3. **Enhancing User Experience:** It aims to simplify the process for users by offering intuitive guidance and insights, making advanced analytics more accessible.

This feature is part of Frontline Systems' broader effort to integrate AI into their suite of decision support tools, making it easier for users to leverage powerful optimization and analytics capabilities.

AI Assist Example

To demonstrate this exciting new feature, open the Product Mix example model by clicking Help – Examples – Optimization Models. While the model worksheet is open, click AI Assist. AI Assist opens and automatically enters the prompt, "Please describe the model and check for correctness."

AI Assist first translates your model to a RASON model, Frontline's modeling language, then diagnoses the model type and examines the model for correct formulation.

Please describe the model and check for correctness

Description of constraints diagnosis

Description of the Model

The provided model is an optimization type, specifically focused on determining the optimal mix of products a company should manufacture in order to maximize gross profit. Here's a detailed breakdown:

Model name & description

General Information

- **Model Name**: ProductMix (Opt)
- **Model Description**: Not provided, likely a version of a Product Mix optimization problem.
- **Worksheets**: The model uses two worksheets named "Example 1" and "Example 2", but it is set to use "Example 1" as the active sheet.

Description of objective

Engine settings

Engine Settings

- **Engine**: LP/Quadratic (Linear Programming/Quadratic)
- **Scaling**: -1 (Automatic scaling is disabled)

Variables

Variables

- **Alias Name**: Number_to_build
- **Initial Value**: 0
- **Lower Bound**: 0
- **Final Value**: []

These variables represent the number of different products to be manufactured.

Description of data used in model.

Data

- **Production Coefficients**: `c18:e22` represents the production coefficients matrix, which is a 2D array of the resources consumed by each product. The values are:
- `[[1, 1, 0], [1, 0, 0], [2, 2, 1], [1, 1, 0], [2, 1, 1]]`
- **Inventory/Resources Available**: `h18:h22` represents the available quantity of each part/resource. The values are:
- `[450, 250, 800, 450, 600]`, referred to as "Number_available".
- **Profit per Unit**: `c24:e24` indicates the profit per unit of each product. The values are:
- `[[75, 50, 35]]`

Constraints

- **Alias Name**: Number_used
- **Formula**: For each part/resource, the total used is computed via the SUMPRODUCT of production coefficients and the variables (Number_to_build).
- **Upper Bound**: "Number_available"

The constraints ensure that the total number of each part/resource used to produce the products does not exceed the available quantities.

```
#### Objective
- **Alias Name**: Total_profit
- **Formula**: The objective function is to maximize
the total profit calculated as the SUMPRODUCT of the
profit per unit and the number of units to be built.
- **Type**: Max (Maximize the total profit)
- **Final Value**: []
```

```
### Checking for Correctness
The model appears to follow the structure of a
standard linear programming problem for product mix
optimization. However, to ensure correctness, here
are some specific checks:
```

```
1. **Data Integrity**: The data values provided for
parts, inventory, and profit are realistic and
consistent with typical optimization models.
2. **Variable Definitions**: Variables for the number
of products to build have correct initial values and
bounds.
```

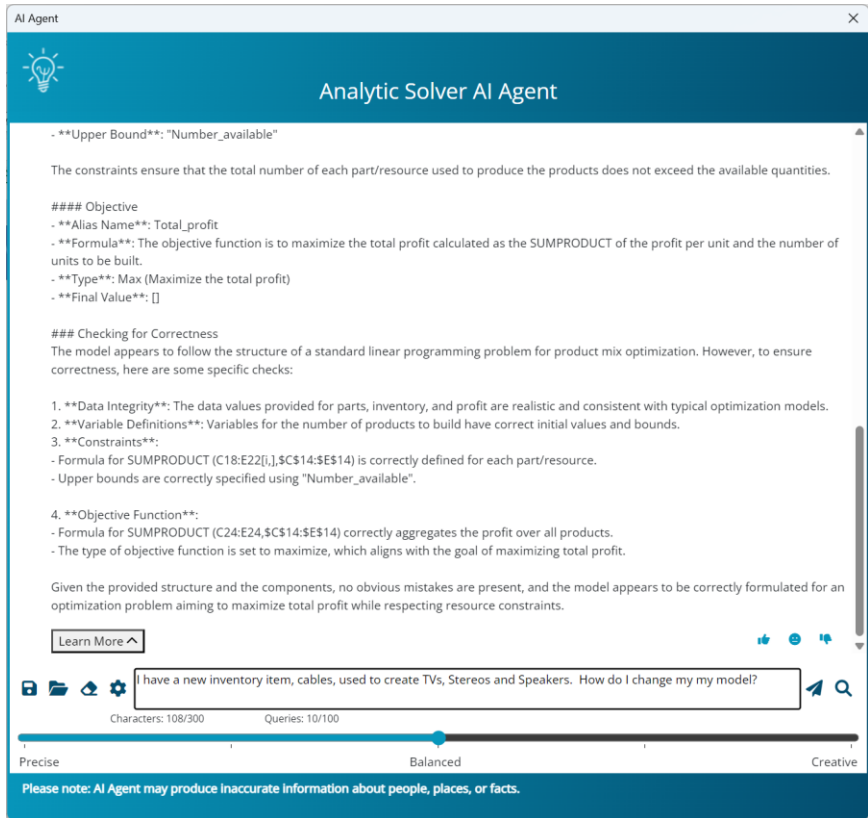
Conclusions
of model
analysis.

```
3. **Constraints**:
- Formula for SUMPRODUCT (C18:E22[i,],$C$14:$E$14) is
correctly defined for each part/resource.
- Upper bounds are correctly specified using
"Number_available".
```

```
4. **Objective Function**:
- Formula for SUMPRODUCT (C24:E24,$C$14:$E$14)
correctly aggregates the profit over all products.
- The type of objective function is set to maximize,
which aligns with the goal of maximizing total
profit.
```

Given the provided structure and the components, no obvious mistakes are present, and the model appears to be correctly formulated for an optimization problem aiming to maximize total profit while respecting resource constraints.

Assume that you have a new part that must be incorporated into the production of all products, but you aren't sure how to change the model formulation to reflect this new part. Simply ask AI Assist!



Notice that the prompt, "I have a new inventory item, cables, used to create TV, Stereos and Speakers. How do I change my model?", gives specific information. The more detailed the context provided in the prompt, the better the response will be. Click OK to view the response.

The response includes a step-by-step guide on all the changes that must be incorporated into your model in order to add a new part to the production requirements.

1. A new inventory item must be added.
2. Data containing the number of cables to be used in the product of each product must be updated.
3. A new constraint must be added to ensure that the inventory is not exceeded.

To include a new inventory item like cables in your optimization model for manufacturing TVs, stereos, and speakers, you'll need to update several components of the model. Here's a step-by-step guide on how to modify the model:

Step-by-Step Guide:

1. **Add New Inventory Item:**
 - In your parts inventory section, add a new row for cables.
 - Specify the initial amount of cables available in the inventory.

Updated Inventory Table:

```

\plaintext
Part Name | Inventory | No. Used - TV | No. Used - Stereo |
No. Used - Speakers
-----
Chassis | 450 | 1 | 1 | 0
Screen | 250 | 1 | 0 | 0
Speaker | 800 | 2 | 2 | 1
Power Supply | 450 | 1 | 1 | 0
Electronics | 600 | 2 | 1 | 1
**Cables** | [Specify] | [Specify] | [Specify] | [Specify]
\`

```

2. **Specify Part Requirement by Product:**
 - Determine how many cables are used for each product (TVs, stereos, and speakers).

- Fill in these values under the "No. Used" columns for TV, Stereo, and Speakers.

Example (Assuming the usage is):

```
``plaintext
Cables | 600 | 2 | 1 | 1
````
```

3. **Update the Constraints:**

- The constraints need to be updated to include the cable usage and inventory.

Updated Constraints in Pseudocode:

```
``plaintext
Chassis_Required <= Chassis_Inventory
Screen_Required <= Screen_Inventory
Speaker_Required <= Speaker_Inventory
PowerSupply_Required <= PowerSupply_Inventory
Electronics_Required <= Electronics_Inventory
Cables_Required <= Cables_Inventory
````
```

In Excel, add a new constraint for the cables in the constraints section:

```
``plaintext
Used_Cables = SUM(No. Used - TV * TVs_Built + No. Used - Stereo * Stereos_Built + No. Used - Speakers * Speakers_Built) <= Available_Cables
````
```

4. **Adjust Objective Function if Necessary:**

- Generally, the addition of a new inventory item does not change your objective function (e.g., maximizing profit), but ensure your objective function is calculated based on the products built.

5. **Update the Excel Model (or other LP modeling tool):**

- Add the new "Cables" row in the appropriate section for parts.  
- Define the usage of cables for each product.  
- Ensure the constraint formulas in Excel or elsewhere include the new cables resource.

6. **Formulate the Model:**

- Double-check if all parts and their respective constraints are properly entered.  
- Ensure the Solver (or other optimization tool) is configured correctly to consider the new constraints.

### Updating the Model Example in Excel:

Assuming you are using Excel, below is a simplified example:

1. **New Inventory Row:**

- Add to the Inventory section, say cell `B6` for cables inventory:

```
``plaintext
Cables | 600
````
```

2. **Usage per Product Row:**

- TV, Stereo, and Speakers in respective locations (e.g., row 12-14 for TV, row 15-17 for Stereo, etc.):

```
``plaintext
Cables Usage | 2 | 1 | 1
````
```

3. **Constraints Update in Solver:**

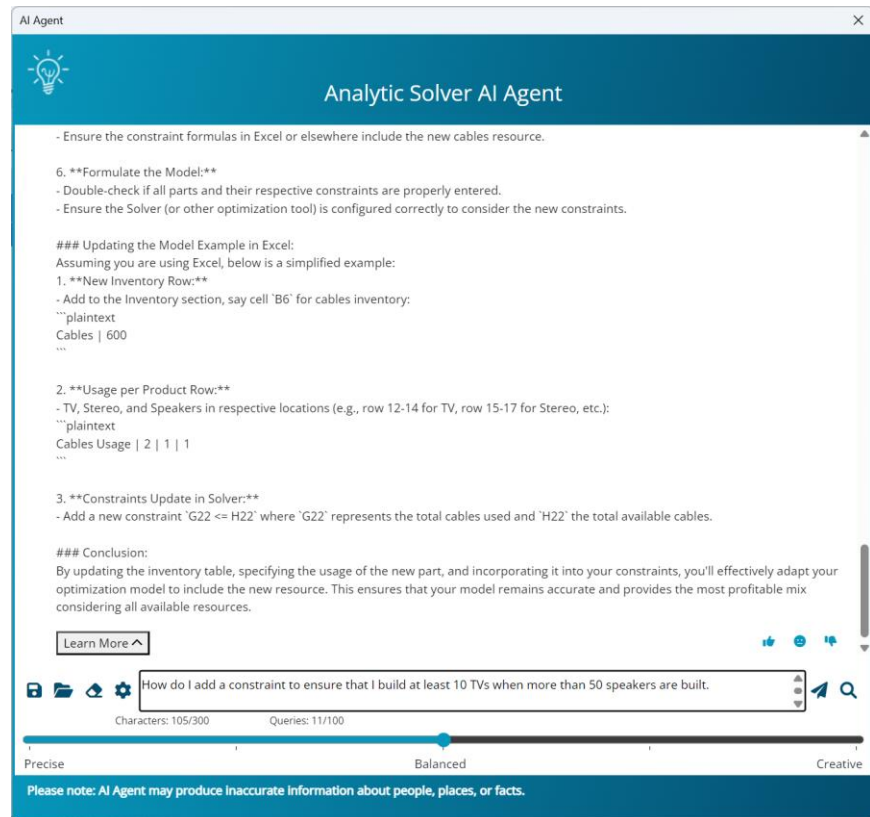
- Add a new constraint `G22 <= H22` where `G22` represents the total cables used and `H22` the total available cables.

### Conclusion:

By updating the inventory table, specifying the usage of the new part, and incorporating it into your constraints, you'll effectively adapt your optimization model to include the new resource. This ensures that your model remains accurate and provides the most profitable mix considering all available resources.

Assume now that for every 50 speakers manufactured, 10 TVs must also be produced. These kinds of real-world "IF Statements" frequently appear in user models, and their conditions are often not intuitive to implement in a linear fashion. Enter AI Assist!

Enter the prompt, "How do I add a constraint to ensure that I build at least 10 TVs when more than 50 speakers are built."



#### AI Assist Result:

The response includes a step-by-step guide on all the changes that must be incorporated into the model in order to add the new constraint.

4. Introduce a new binary variable.
5. Add a new constraint that ensures that if the number of speakers built is greater than 50, then the binary variable is 1 and 0 otherwise.

To add a constraint ensuring that you build at least 10 TVs ( $x_1$ ) when more than 50 speakers ( $x_3$ ) are built in your linear programming model, you can use a combination of binary variables and a large constant ( $M$ ). This is a common approach for creating conditional constraints in linear programming models.

Here's the step-by-step method to add this constraint:

1. **Introduce a Binary Variable ( $y$ ):**
  - Let  $y$  be a binary variable that indicates whether more than 50 speakers are being built.  $y$  will be 1 if more than 50 speakers are built and 0 otherwise.
2. **Add a Big M Constraint:**
  - Use a sufficiently large constant  $M$  to enforce the conditional logic.  $M$  should be larger than the possible number of speakers that could be built (e.g., 100 here).
3. **Formulate the Conditional Constraint:**
  - The constraint that ensures at least 10 TVs are built

when more than 50 speakers are built can be written as:

```
\[
x1 \geq 10y
\]
- \((x3 - 50 \leq My \))
- This ensures that \((y \)) must be 1 when more than 50
speakers are built, triggering the \((x1 \geq 10 \))
constraint.
```

4. **Implement in Your Solver Model:**  
 - Add the binary variable \(( y \)) and the constraints in your Solver tool.

Here is how you can structure these constraints in general terms:

1. **Binary Variable Constraint:**  
 - Add a binary variable constraint:

```
\[
y \in \{0, 1\}
\]
```

2. **Relationship between \(( y \)) and \(( x3 \)):**

```
\[
x3 - 50 \leq My
\]
```

3. **Ensuring At Least 10 TVs if \(( y = 1 \)):**

```
\[
x1 \geq 10y
\]
```

### In Excel or Solver Model Interface

Assuming you are using Excel Solver, follow these steps:

1. **Define a Binary Variable:**  
 - Assume you add it in cell G9, and it will be called \(( y \)).

2. **Add the Constraints:**

a. For \(( y \)):

- **Cell:** G9  
 - **Constraint:** Bin (Binary)

b. For the relationship between \(( y \)) and \(( x3 \)):

- **Formula Cell:** H9 (Formula: `=F9 - 50 - M*G9 <= 0`)

c. For ensuring at least 10 TVs:

- **Formula Cell:** I9 (Formula: `=D9 - 10*G9 >= 0`)

3. **Integrate into Solver:**

- Go to the Solver Parameters dialog box and add the constraints created in steps 2a-c.

### Example Configuration in Excel Solver:

If:

- \(( x1 \)) (Number to Build: TVs) is in cell D9,  
 - \(( x3 \)) (Number to Build: Speakers) is in cell F9,  
 - Binary variable \(( y \)) is in cell G9,  
 - Using a large number \(( M = 100 \))

Constraints added in Solver:

1. **Binary Constraint for \(( y \)) (Cell G9):**

- Set \(( G9 \)) to be a binary variable.

2. **Big M Constraint:**

- \(( H9 \)): \(( =F9 - 50 - 100\*G9 <= 0 \))

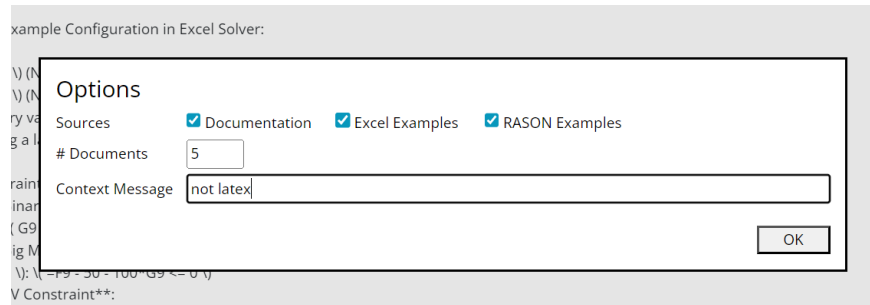
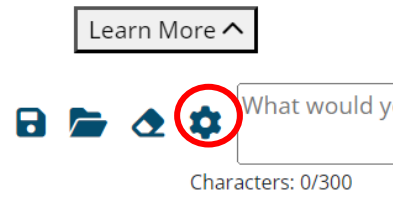
3. **TV Constraint:**

- \(( I9 \)): \(( =D9 - 10\*G9 >= 0 \))

This configuration ensures that your linear programming

model will enforce the requirement of building at least 10 TVs if more than 50 speakers are built.

Note that the syntax in the results is LaTeX. To turn off this syntax, click the gear icon to open the Options dialog, type "not latex" into the Context Message field and then click OK.



Click OK to view the results in a more readable format.

To accommodate a new inventory item, specifically cables, in your existing linear programming model, you will need to follow these steps:

1. **Add the New Inventory Item (Cables) to the Model:**
  - Determine the initial inventory level for cables.
- ...

As with the AI Agent, toggle between Precise, Balanced, and Creative to determine the type of answers returned. If Precise is selected, AI Agent will attempt to be as exact and deterministic as possible while Creative will usually result in more original, uncertain and non-repetitive answers. Use Balanced (the default) for the best of both worlds.

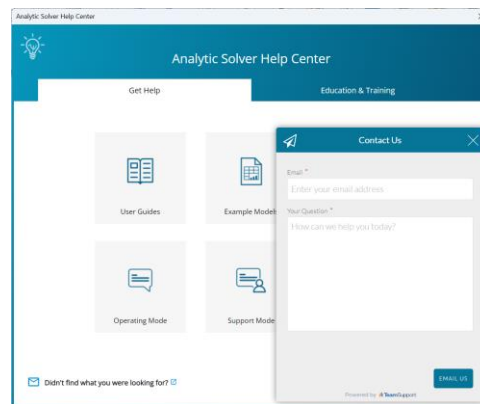
## Maximizing Benefits of AI Agent and AI Assist

- Use “New Topic” frequently if your question does not require chat history. Results will be faster and more precise.
- Be as concise as possible in your question, i.e. “How to use RASON?” and “Give me step-by-step instructions on how to use RASON” will produce different results.
- Analytic Solver Upgrade or Solver Academy licenses do not include support for AI Agent. This button will be disabled.

Note: The AI Agent in Frontline Solvers software uses “generative artificial intelligence” methods, which can at times produce erroneous information about people, places or facts, including incorrect information about Frontline’s software products. Users are advised to take this into account, and use “common sense and human intelligence” when conversing with the AI Agent. Frontline Systems Inc. assumes no responsibility for inaccuracies in the AI Agent’s responses.

## Help Center

Click Help – Help Center to open the Help Center. Click *Support Live Chat*, in the bottom right hand corner, to open a Live Chat window. If you run into any issues when using the software, the best way to get help is to start a Live Chat with our support specialists. This will start a Live Chat during our business hours (or send us a message at other hours), just as if you were to start a Live Chat on [www.solver.com](http://www.solver.com) – but it saves you *and* our tech support rep a lot of time – because the software reports your latest error message, model diagnosis, license issue or other problem, without you having to type anything or explain verbally what’s happened. You’ll see a dialog like this:

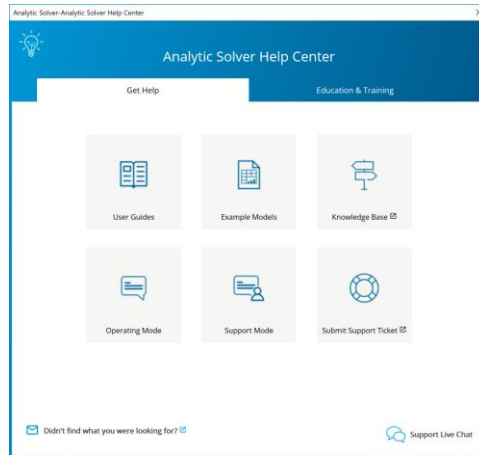


Since the software automatically sends diagnostic information to Tech Support, we can usually identify and resolve the problem faster. (Note: No contents from your actual spreadsheet model is sent, only information such as the number of variables and constraints, last error message, and Excel and Windows version.)

Note: If Support Live Chat is disabled, click the down arrow beneath Help and select *Support Mode – Active Support*.

## Accessing Resources

The Help Center gives you easy access to video demos, User Guides, online Help, example models, and Website support pages to learn how to use our software tools, and build an effective model.



## User Guides

Click the User Guides menu choice to open PDF files of the Analytic Solver Optimization and Simulation User and Reference Guides, Analytic Solver Data Science User or Reference Guides, or our Quick Start Guides.

## Example Models

Clicking this menu item will open the Frontline Solvers Example Models Overview dialog with nearly 120 self-guided example models covering a range of model types and business situations.

## Knowledge Base

Click Knowledge Base to peruse a multitude of online articles related to support and installation issues or to locate articles that will help you to quickly build accurate, efficient optimization, simulation, and data science models.

## Operating Mode

Click Operating Mode to switch between three different levels of help. The Excel formulas and functions you use in your model have a huge impact on how fast it runs and how well it solves. If you learn more about this, you can get better results, but if you don't, your results will be limited. Guided Mode can help you learn.

- Guided Mode prompts you step-by-step when solving, with dialogs.
- Auto-Help Mode shows dialogs or Help only when there's a problem or error condition.
- Expert Mode provides only messages in the Task Pane Output tab. (This mode not supported when using a trial license.)

## Support Mode

Click Support Mode to switch between three different levels of support. No information (cell contents etc.) from your Excel model is ever reported automatically to Frontline Systems, in any of these Support Modes. Only events

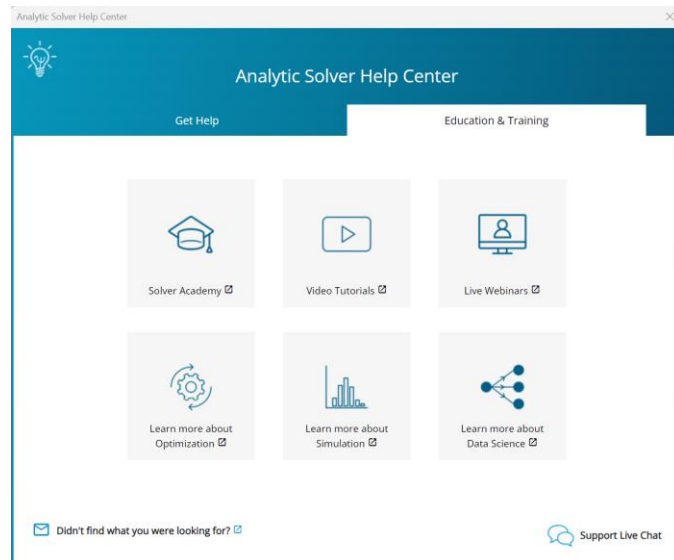
in Frontline's software, such as menu selections, Solver Result messages, or error messages are reported.

- **Active Support** automatically reports events, errors and problems to Frontline Support, receives and displays messages to you from Support, and allows you to start a Live Chat with Support while working in Excel (**Recommended**). Auto-Help Mode shows dialogs or Help only when there's a problem or error condition.
- **Standard Support** automatically reports events, errors and problems anonymously (not associated with you) to Frontline Support, but does not provide a means to receive messages or start a Live Chat with Support.
- **Basic Support** provides no automatic connection to Frontline Support. You will have to contact Frontline Support manually via email, website or phone if you need help.

## Submit a Support Ticket

If you're having installation, technical, or modeling issues, submit a Support Ticket to open an online support request form. Submit your email address and a short, concise description of the issue that you are experiencing. You'll receive a reply from one of Frontline's highly trained Support Specialists within 24 hours, and generally much sooner.

Our technical support service is designed to supplement your own efforts: Getting you over stumbling blocks, pointing out relevant sections of our User Guides or example models, helping you fix a modeling error, or -- in rare cases -- working around an issue with our software (always at our expense).



## Solver Academy

[Solver Academy](#) is Frontline Systems' own learning platform. It's the place where business analysts can gain expertise in advanced analytics: forecasting, data science, text mining, mathematical optimization, simulation and risk analysis, and stochastic optimization.

## Video Tutorials/Live Webinars

Click Video Tutorials to be directed to Frontline's You Tube Channel. Browse videos on how to create an optimization or simulation model or construct a data science or prediction model using Analytic Solver.

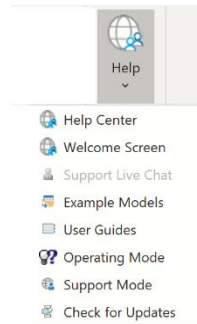
Click Live Webinars to be redirected to [www.solver.com](http://www.solver.com) to join a live or pre-recorded webinar. Topics include *Using Analytic Solver Data Mining to Gain Insights from your Excel Data*, *Overview of Monte Carlo Simulation Applications*, *Applications of Optimization in Excel*, etc.

## Learn more!

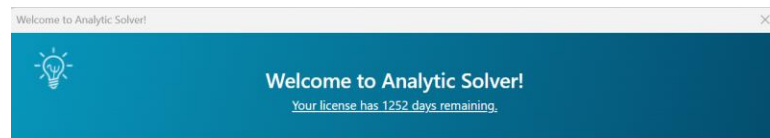
Click any of the three Learn More buttons to learn more about how you can solve large-scale optimization, simulation, and data science models, reduce costs, quantify and mitigate risk, and create forecasting, data science and text science models using Analytic Solver.

## Help Menu

Use the Help Menu to gain short cuts to live chat, example models, documentation, set your operating and support mode preferences, and also to open the Welcome Screen and check for software updates.



Use the Welcome Screen to get help with an existing model, open our example models or watch a quick video on how to get running quickly with Analytic Solver.



To get started, click a box below






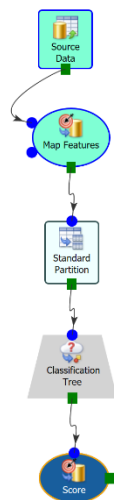
# Creating & Deploying a Workflow

---

## Introduction

The Workflow tab, released in version 2018, allows the combination of all available data science techniques into an all-inclusive workflow, or workflows. Once the workflow is created, either manually or simply by recording your actions, you can initiate the start of the pipeline, or pipelines, by clicking the  button. Note: This is the only tab supported in the Data Science Cloud app.

Analytic Solver Comprehensive allows you to export your existing workflow to Frontline's Rason Cloud Services where it can be deployed to a website without the need to write code! Continue reading to find out how.



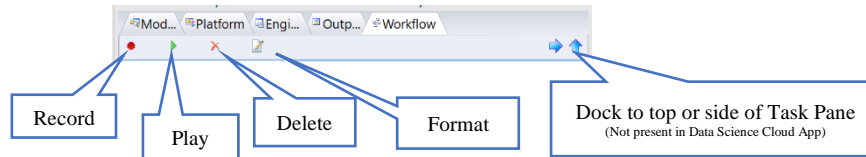
---

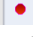
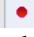
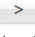
## Creating a Workflow

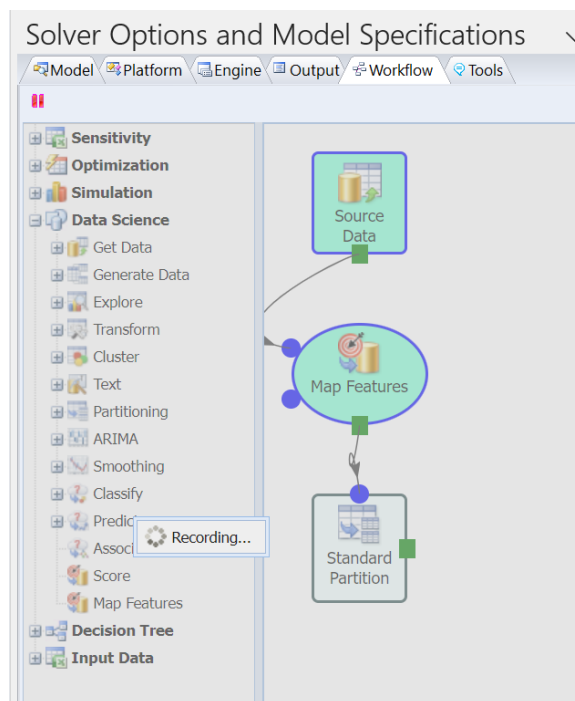
Workflows can be created in two ways: 1. By recording your actions performed on the Data Science ribbon after the Record button is pressed on the Workflow tab or 2. By manually dragging the Data Science nodes from the left of the Workflow tab onto the Workflow window. The next two sections, Recording a Workflow and Manually Creating a Workflow, contain examples to illustrate each method.

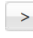
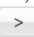
### ***Recording a Workflow***

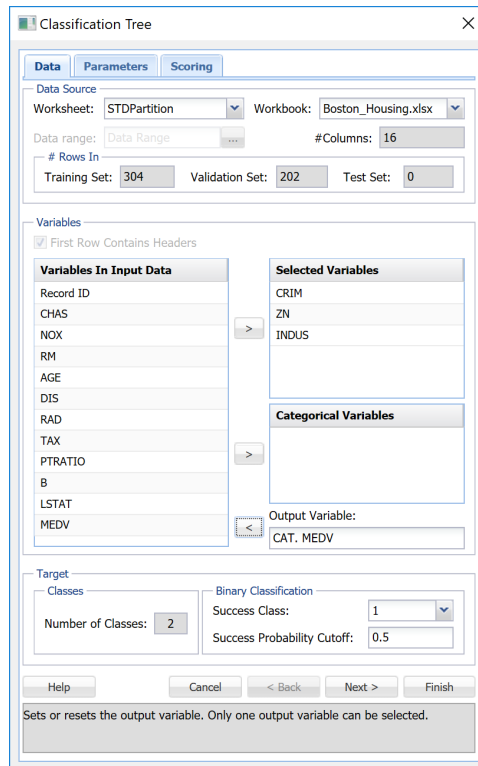
At the top of the Workflow tab, you'll find 6 icons. Each one performs a specific function.



**Press  to record a new workflow.** On the Data Science ribbon, click Help – Example Models – Forecasting/Data Science Examples, then click the Boston\_Housing dataset hyperlink to open the Boston Housing example dataset. Click a cell on the Data worksheet, then click  to start the Workflow recorder. Click Data Science – Partition – Standard Partition to open the Standard Data Partition dialog. While holding down the shift key, select all variables under *Variables In Input Data*, then click  to move them to *Selected Variables*. Leave all options at their defaults, then click OK to run the partition. You'll see several icons added to the Workflow tab.

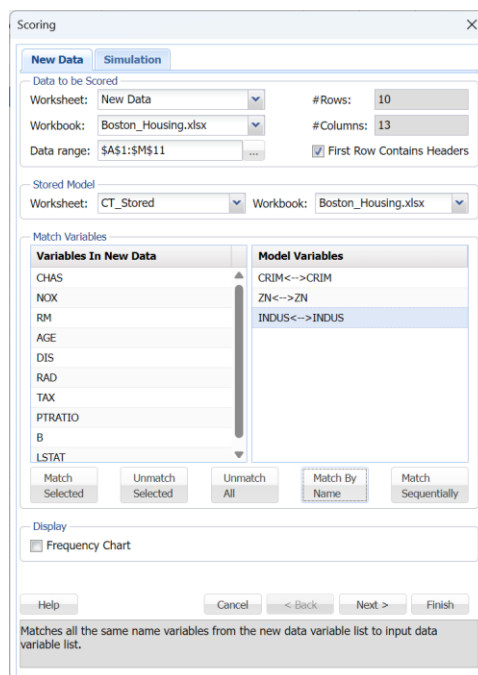


Click Data Science -- Classify -- Classification Tree to open the Classification Tree dialog. Under Data Source, click the down arrow next to Worksheet and select STDPartition. While holding down the SHIFT key, select CRIM, ZN, and INDUS, under *Variables In Input Data*, then click  to move these three variables to *Selected Variables*. Next, select the CAT.MEDV variable, under *Variables In Input Data*, and click  to choose this variable as the *Output Variable*.

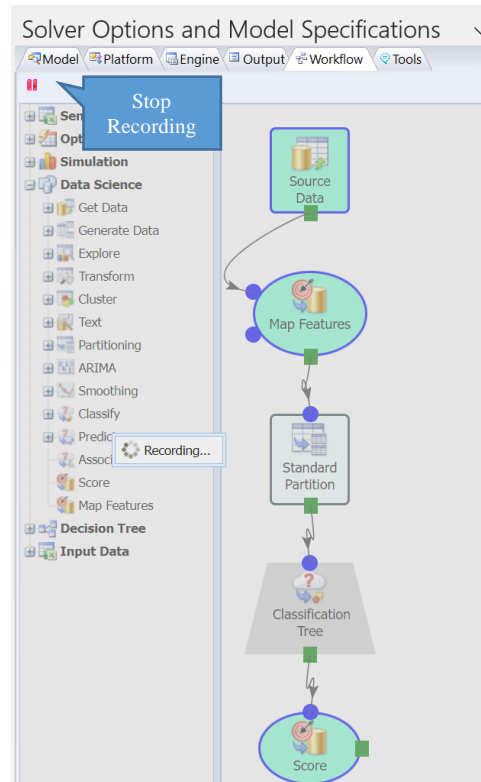


Leave all remaining options at their defaults, then click Finish to run the Classification Tree algorithm.

Finally, click the Score button on the Data Science ribbon to open the Scoring tab. Under *Data to be Scored*, click the down arrow next to Worksheet and select New Data. The *Variables In New Data* field will populate with new variables. Afterwards click *Match By Name* to match the variables listed under *Variables in New Data* with the variables listed under *Model Variables*.




Click *Finish* to score the new data then click the Stop Recording button to stop the Workflow recorder.





The completed workflow is displayed in the window. Use Excel’s File – Save to save the workflow to the workbook.

There is no limit to the number of workflows that may appear in the Workflow tab. Each new workflow will be added to the right of the existing workflow(s).

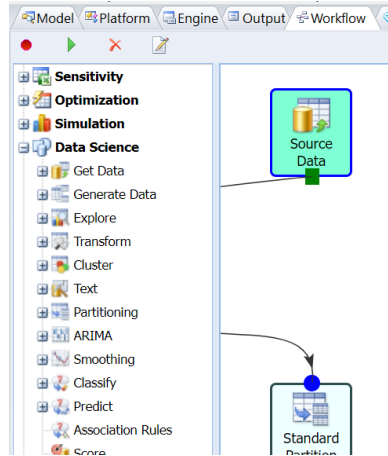
## Workflow Tab Options

Press  to execute the workflow or workflows. If no nodes in the workflow are selected, all workflows present in the Workflow window will be executed. If a node in the flow is selected, and the Execute Workflow button is pressed, the workflow will be run up to and including the selected node.



Press  to delete the contents of the Workflow window. If no nodes in the workflow are selected, all nodes in the Workflow window will be deleted. If a node is selected in the Workflow window, pressing the Delete Workflow icon will delete only the selected node.

Press  to format the Workflow into a more readable form. If the Task Pane is docked to the top of the Excel workbook, the Workflow will be formatted horizontally. If the Task Pane is docked to the side of the Excel workbook, the Workflow will be formatted vertically.

Note: It is possible to move a node “behind” the task pane as shown in the screenshot below.



To regain access to the node, click the Format Workflow button.

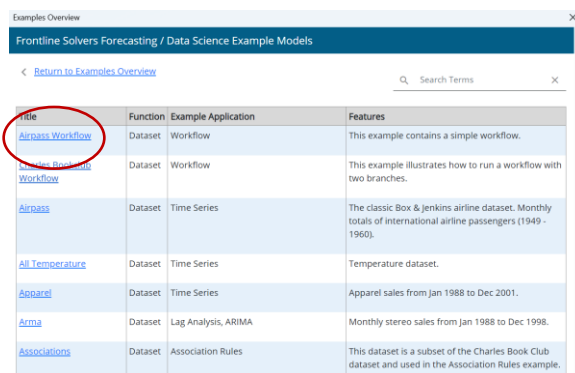
Press  to dock the Solver Task Pane to the right of the workbook. Press  to dock the Solver Task Pane to the top of the workbook.

## Deploying Your Workflow

Starting with Analytic Solver Comprehensive V2020, when you click Deploy Model – Cloud Service, your model (or a copy of it) is saved to the Azure cloud. (Note: Analytic Solver Data Science does not support this new feature. Trial licenses of Analytic Solver Comprehensive *do* support this new feature.)

If you use File Save to store your workbook in OneDrive or OneDrive for Business, we use that, which is by far the best option. Afterwards, any web or mobile application can easily run your model, via simple, standard REST API calls, often made from JavaScript on a web page, of from C# or Java code on a Web server.

Click Help – Example Models – Forecasting / Data Science Examples, then click the Airpass Workflow link. The Airpass Workflow Example will open.

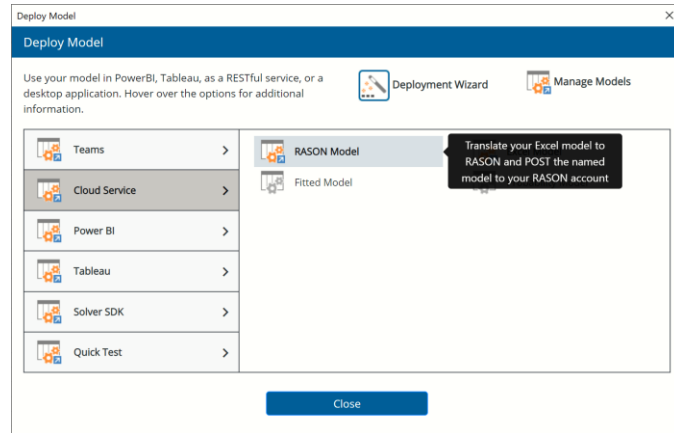


Click the Analytic Solver tab to move to the Analytic Solver ribbon.



# Posting Workflow to RASON Cloud Services

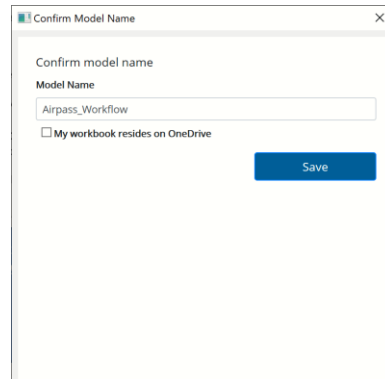
Click Deploy Model– Cloud Service – RASON Model.



In the COM Addin, it's easy for us to check if the user has an optimization, simulation, or flow model, and we can give that proper warning/error message.

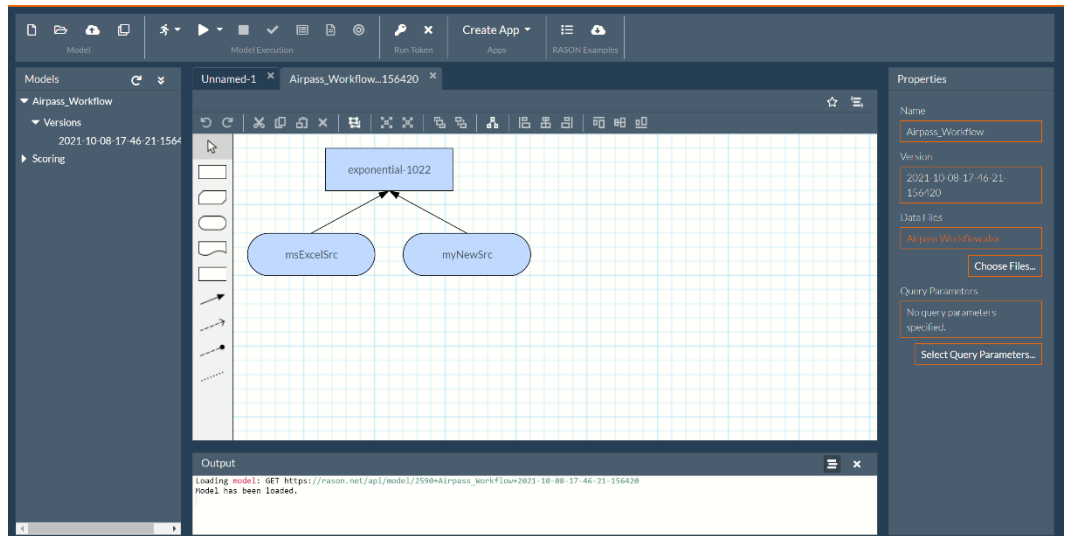
Note: *For a multi-stage data science flow, Cloud service – RASON model is the only available Deploy Model options in either Analytic Solver Com Addin or Analytic Solver Cloud app.* Users of the Analytic Solver COM Addin will receive an error if they attempt to use Power BI, Tableau, etc. However, users of the Analytic Solver Cloud app will not.

On the next screen, click Save to accept the default Model Name. (You can also provide a more meaningful name here if wanted. If your workbook resides on your OneDrive account, click the checkbox. )

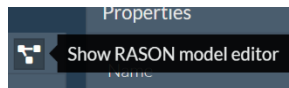


Note that data science models that include a weight variable or a partition variable may not be translated to Rason.

Immediately, a browser opens to the Editor tab on [www.Rason.com](http://www.Rason.com). From here you can view your workflow translated into Frontline's Rason modeling language. Notice that the Airpass\_Workflow is listed under Models on the right and the workflow appears in graphical form in the grid.



To see the model translated into Rason simply click the Show RASON model editor icon in the top left of the grid.



The Rason model editor displays the Airpass Workflow example translated into the Rason modeling language.

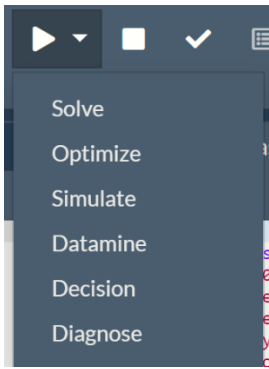
```

1- {
2- "flow": "Airpass_Workflow",
3- "exponential-1022": {
4- "datasources": {
5- "msExcelSrc": {
6- "type": "excel",
7- "connection": "Airpass Workflow.xlsx",
8- "content": "time-series",
9- "selection": "Data!A2:B145",
10- "headerExists": true
11- },
12- "myNewSrc": {
13- "type": "excel",
14- "connection": "Airpass Workflow.xlsx",
15- "selection": "Data!A2:B145",
16- "headerExists": true
17- }
18- },
19- "datasets": {
20- "myData": {
21- "binding": "msExcelSrc"
22- },
23- "myNewData": {
24- "binding": "myNewSrc",
25- "selectedCols": ["Passengers"]
26- }
27- },
28- "modelName": "exponential-1022",
29- "modelType": "datamining",
30- "estimator": {
31- "exponential-1022": {
32- "type": "timeSeries",
33- "algorithm": "exponential",
34- "parameters": {
35- "optimize": true
36- }
37- }
38- }
39- }
40- }

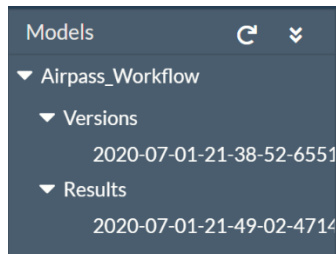
```

From here you can edit your model directly in Rason or if you'd rather not interact with Rason code at all, simply make the change to the Excel model and then re-deploy the edited workflow. There's no need to learn a new language!

To solve this model on the Editor tab, simply click the down arrow next to the "play" button and click Solve. (Only the Solve endpoint supports the solving of workflows.)

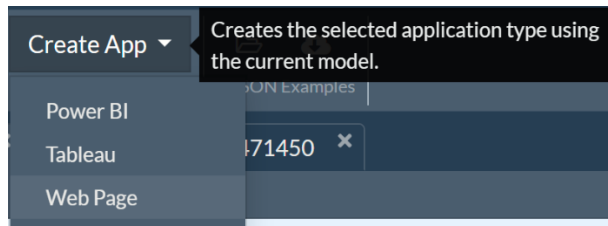


After the workflow is executed, you can find the results listed under Results.

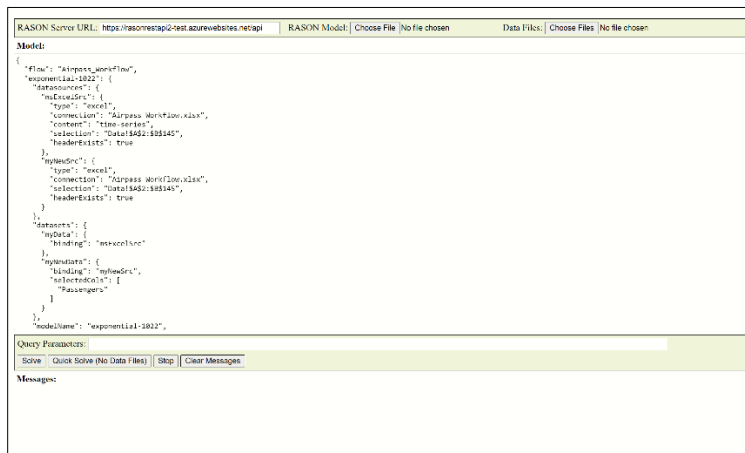


## Deploying Your Model

Click Create App – Web Page on the Editor tab ribbon to deploy the Airpass Workflow to a Web page application.



The file, RasonScript.html, will be downloaded locally. Double click the file to open the Web application. Notice that the translated workflow appears in the Model window.

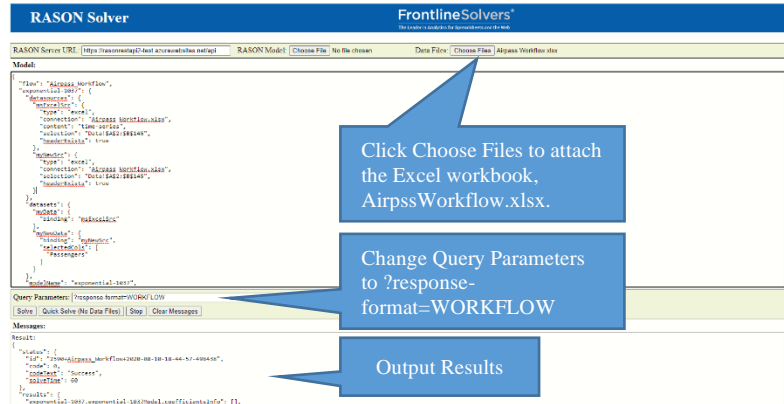


To solve the workflow:



1. Attach the Excel file, Airpass Workflow.xlsx. Click "Choose Files" at the top of the Web app, then browse to C:\ProgramData\Frontline Systems\Datasets and select Airpass Workflow.xlsx.
2. Since we are solving a workflow, we must first change the query parameter from ?response-format=STANDALONE to response-format=WORKFLOW.
3. Now click the Solve button to solve the workflow through the Web application.

Note: Quick Solve does not support external files so it is unable to solve this translated workflow.

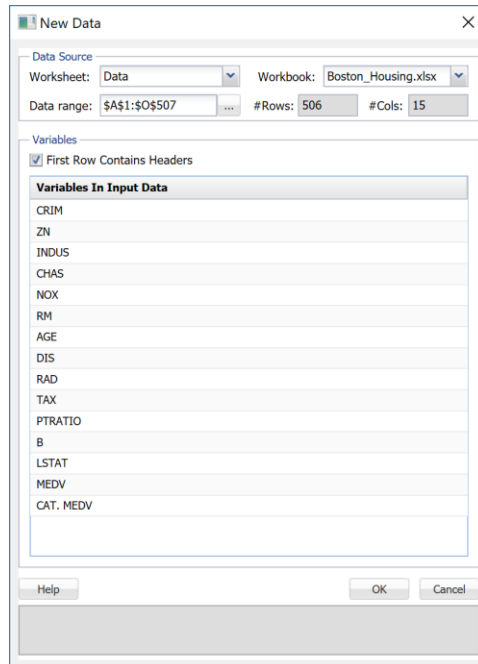


Whether you decide to keep your model in Excel or translate it to RASON, the *important thing* is that your workflow (or model) can be executed **whenever or wherever** it's needed – on the factory floor, on a salesperson's laptop or smartphone or in call center custom application. And Rason can get the updated data your model needs directly from operational business systems. This is amazingly easy if you are using Power BI, Power Apps, Power Automate or Dynamics 365.

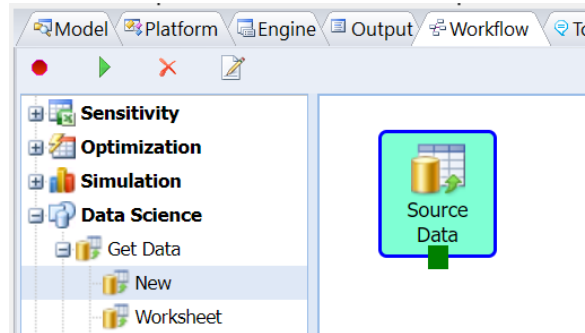
For more information on this new feature, including information on the Deployment Wizard and how to Manage your deployed models within Excel, see the Analytic Solver User Guide chapter, Deploying Your Model.

## Manually Creating a Workflow

It's also possible to manually create a workflow or workflows. Go back to Excel, open 'Boston Housing.xlsx' workbook and navigate to 'Workflow' tab. Click the + in front of Get Data, then drag "New" into the workflow window. The New Data dialog will open. Under Data Source, click the down arrow next to Worksheet and select Data.

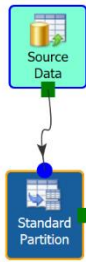


Click OK to close the dialog. The Source Data icon will appear in the Workflow window. It's important to note that a workflow does not need to start with a Source Data node. In fact, any node dragged to the workflow window can stand alone and will be executed when the Execute Solver button is pressed.



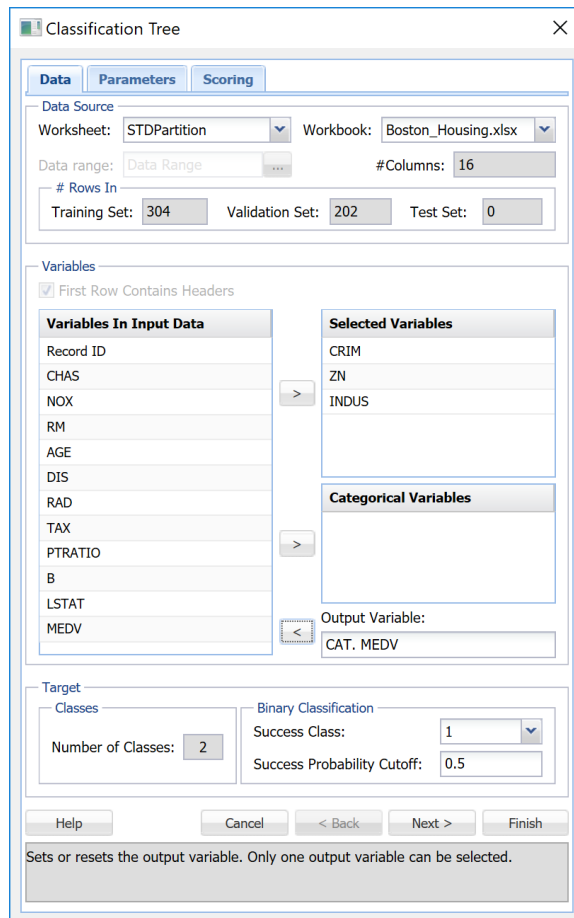
Afterwards, expand Partitioning (by clicking the +), then drag Standard Partition to the Workflow window. When the Standard Partition dialog opens, select all variables under *Variables In Input Data* then click **>** to move all variables to *Selected Variables*. Leave all options at their defaults and click OK. Immediately, the workflow updates in the Workflow window displaying the connection.

Note: Notice that when creating a workflow manually, Map Features is not required. However, if you'd like to provide a new data source for an existing flow, you must use Map Features to specify the input and output variables to be used in the flow. See below for more information.



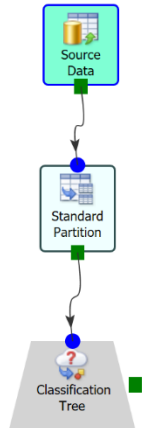
Expand Classify, then drag Classification Tree to the window.

- Under Data Source, click the down arrow next to Worksheet and select STDPartition.
- Select the first three variables under *Variables In Input Data* (CRIM, ZN and INDUS) and click **>** to move them to *Selected Variables*.
- Select CAT.MEDV as the *Output Variable*.



Then click Finish. For more information on the Classification Tree feature in Analytic Solver Data Science, see the Data Science Reference Guide.

The workflow will update immediately by connecting the Standard Partition icon to the Classification Tree icon.

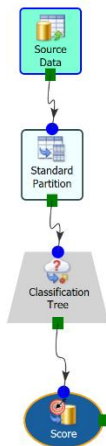



Note: Automatic connections are made based on the current effective state of the workflow and sheets in the workbook. For example, if you were to run this workflow, a new STDPartition1 would be created. As a result, connecting a new classification or regression to the STDPartition sheet would not be allowed.

On the Workflow tab, drag Score into the Workflow window. The Scoring dialog will open. Click the down arrow next to Worksheet under *Data to be Scored* and select *New Data*. Then click the down arrow for *Worksheet* under *Stored Model* and select *CT\_Stored*. Click *Match By Name* to match the variables in the New Data worksheet with the selected Model Variables, then click OK.

Note: When a node is dragged onto the workspace, the node's dialog opens. When Finish or OK is pressed on the open dialog, the node is executed. If this dialog is reopened and either minor or no changes are made, the node will not be executed again when Finish or OK is clicked to close the dialog. However, if major changes are made, and children (additional nodes) are attached to this node, the node will be executed again.

Immediately, the workflow window is updated.



Once all nodes are connected, click  to run the workflow.

## ***Saving a Workflow***

A workflow is automatically saved to the workbook when one of three events is executed.

1. When an Analytic Solver Data Science dialog is closed.
2. After the workflow is formatted by pressing the Format Workflow button on the Workflow tab.
3. After the Stop Recording button is pressed on the Workflow tab.

Otherwise, use Excel's File – Save menu item to save the workflow to the workbook.

## ***Note on Sampling/Summarizing Big Data***

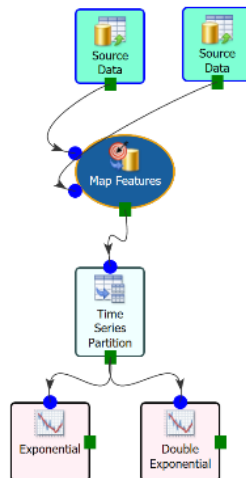
Currently, workflow methods that perform asynchronous activities, such as Sampling and Summarizing Big Data, will first be converted to synchronous methods when included in a workflow. Running big data synchronously (i.e. clicking the RUN button on the Sampling or Summarizing Big Data dialog) is supported without any required internal conversion.

## **Making/Breaking a Connection**

To make a connection between two existing nodes, simply connect the green square on the first node with the blue circle on the second node. To break a connection, simply click the arrow connection in the blue node of the second node and move it back to the green square on the first node.

## **Running a Workflow with a New Dataset**

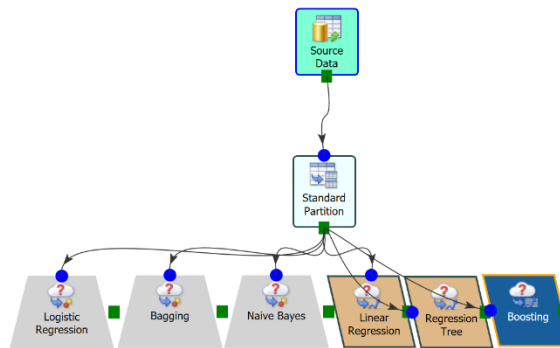
When there is an existing workflow, you may drag and drop a new dataset (a new Get Data node) onto the Workflow window and connect the new Source Data icon to the existing “Map Features” node. The result is two connections, one from each Source Data icon, to the Map Features icon. Connecting to, or double-clicking the “Map Features” node will pop up the existing Score button dialog, enabling column names to be matched. When the workflow is executed, the newer data source will be used.



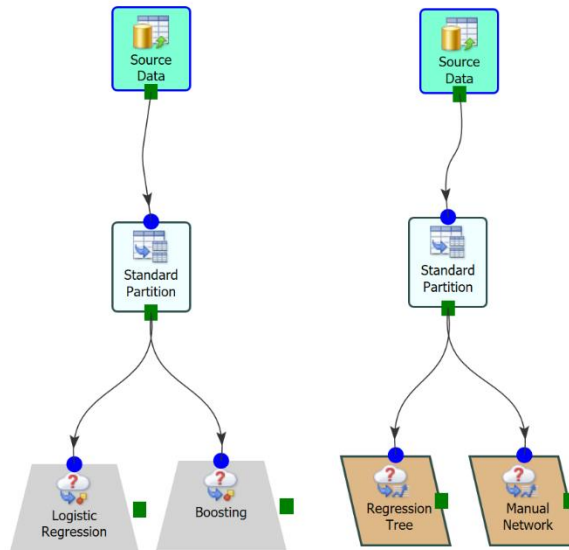
## Multiple Workflows

As discussed earlier, there is no limit to the number of workflows appearing in the workflow window. To run several different classification or regression algorithms at the same time, simply drag the desired method into the workflow and connect the method to the flow.

In the example workflow below, three different classification techniques, Logistic Regression, Bagging and Naïve Bayes, and three different regression methods, Linear Regression, Regression Tree and Boosting, use the same standard partitions in the workflow. When this flow is executed, all six models will be created and used to score new data. Note: Although each classification and regression technique uses the same standard partition and data source, each model may contain different selected variables.



Two workflows are present below. Once the Execute Workflow button is pressed, all workflows will be executed at the same time.

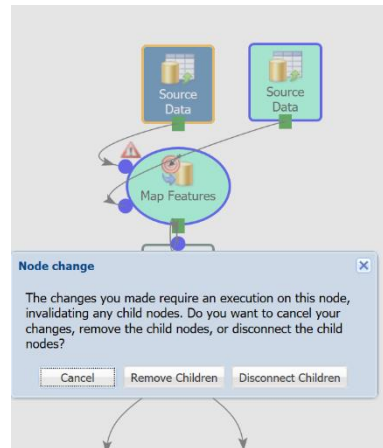


## Changing Options Settings

To change the option settings for any of the nodes in the workflow, simply double click the desired node to bring up the task dialog, make the desired option changes, then click OK or Finish.

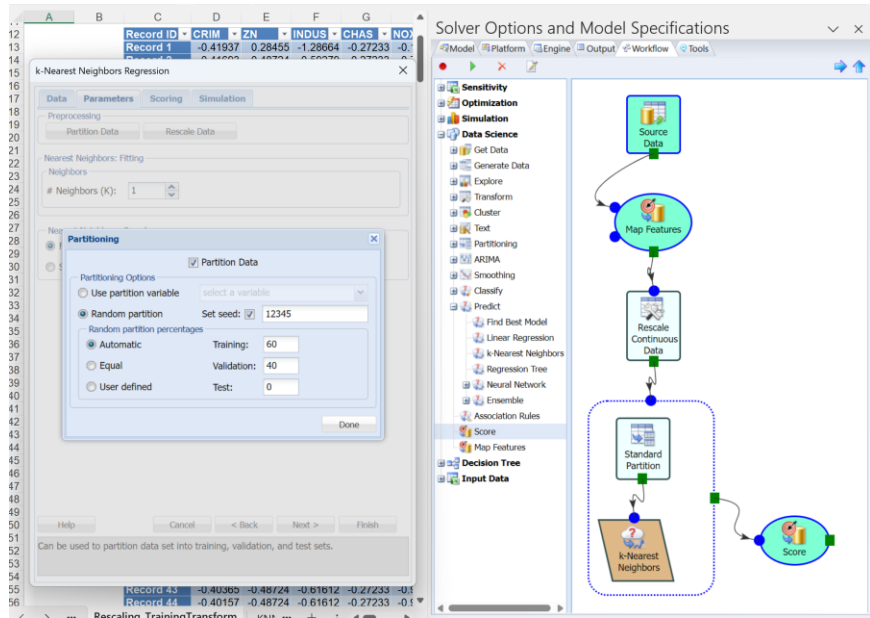
**Important Note:** If you make changes to the Selected Variables or the Input Sheet for any existing node, a message with three options will appear asking if you would like to 1. Cancel, 2. Remove Children or 3. Disconnect Children.

Click Cancel to return to the Workflow tab without any changes being made. Click Remove Children to delete all nodes beneath the selected node or Disconnect Children to disconnect the selected node from the rest of the workflow.



## Workflow Groups

Analytic Solver Data Science allows you to partition, rescale and score a dataset “on-the-fly”. This means that instead of performing separate steps for partitioning, rescaling or scoring, users can perform all these actions on the same set of dialogs used to create the classification or regression model. When partitioning, rescaling or scoring “on the fly” during the recording of a workflow, a “group” will be created. This group is treated as one node and is denoted with a dotted line surrounding all nodes in the group. Clicking any of the nodes in the group will bring up the classification or regression method dialogs.



To create a group when creating a workflow manually, simply select the partitioning, rescaling or scoring options once the classification or regression method dialog opens.

To remove or change a group, double click any of the nodes in the group to open the classification or regression method dialog and disable the desired feature or features.



# Deploying a Fitted Model to RASON, Power BI or Tableau

---

## Introduction

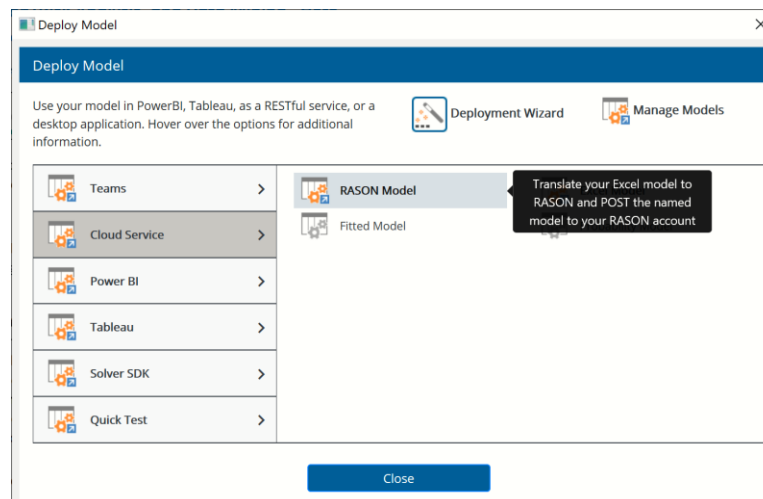
When using Analytic Solver you have the ability to not only create, design and solve your optimization, simulation, simulation optimization, stochastic optimization and data science/forecasting models in Microsoft Excel, but you also have the ability to embed your model into your own application and solve it on the Web or within Power BI or Tableau by calling our RASON Server.

RASON (which stands for Restful Analytic Solver Object Notation) is Frontline Systems' modeling language embedded in JSON and a REST API that makes it easy to create, test and deploy analytic models using optimization, simulation and data science, in web and mobile applications.

Users of Frontline's Excel Solvers will find that 1. it's exceedingly easy to translate Excel models into RASON models 2. that your knowledge of Excel formulas and functions is immediately usable and 3. RASON models can be more flexibly "bound" to data from a variety of sources.

When you click the RASON Deployment icon located on the Analytic Solver ribbon, you'll notice that your optimization, simulation or forecasting/data science model can be deployed to the Cloud, Power BI, Tableau, or your own application on the Web.

Note: In order to use the RASON Deployment functionality, users of Analytic Solver Data Science must also have a license for Analytic Solver Upgrade, Decision, Optimization, Simulation or Comprehensive. For more information on fitted models, please see the Scoring chapter within this guide.



Web App Developers will be able to immediately find how exceptionally easy it is to embed RASON models as JSON and solve them using Frontline's RASON server, which exposes a simple REST API that's scalable to handle very large,

compute-intensive analytic models. Months of work, that would have previously been required, have been reduced to a single command button click!

---

## RASON Deployment Menu

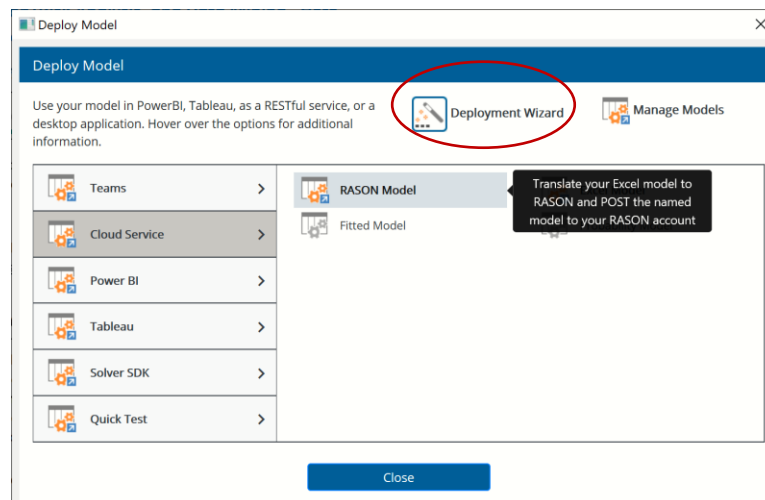
The purpose of RASON Deployment is to enable your model to be used wherever/whenever it is needed – on the factory floor, on a salesperson's laptop or smart phone or in a call center custom application. On your own, without web development or IT help, you can get your model working as a cloud service. Once you've shown that it works, your developers or IT will find they can everything they'd like – customize code, connect to operational databases, apply security and governance best practices.

Rason makes your model accessible 24 x 7 via simply REST API and JSON responses – the common standard for web and mobile apps. Even better, it makes your model results available via OData – the standard widely supported by Microsoft apps including Power Platform. And Rason has built-in facilities to get the updated data your model needs directly from operational business systems. This is amazingly easy if you are using Power BI, Power Apps, Power Automate, or Dynamics 365.

Rason contains the entire Excel formula language as a subset, including virtually all of Excel's built-in functions. That makes translation possible and it allows you to see your Excel formulas embedded in Rason's syntax. To work well in web and mobile apps, Rason is embedded in JSON (JavaScript Object Notation) which allows it to work directly in Power BI, Tableau or a page on any website.

With RASON, you can create more flexible models, by working with Rason arrays and tables instead of fixed-size cell ranges. When products, regions, time periods and other "dimensions" change, your model can automatically adjust and keep running – without any extra work on your part.

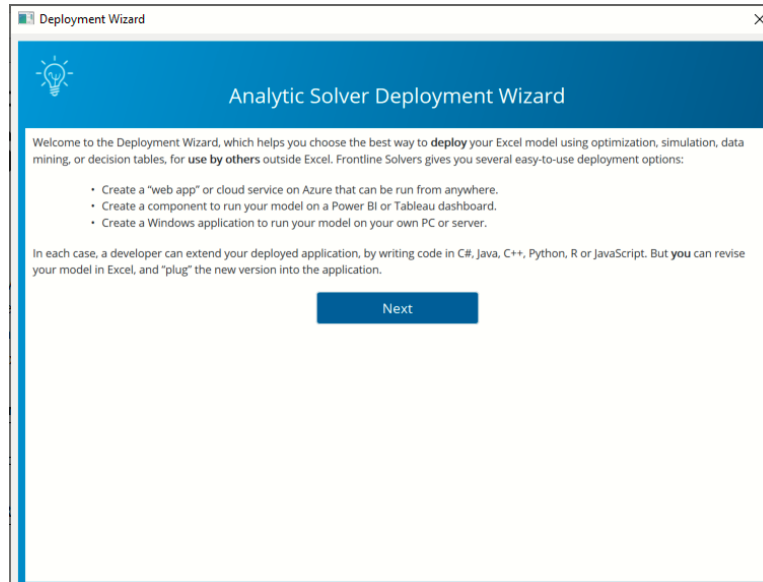
When you first click the RASON Deployment icon on the Analytic Solver ribbon, the new Rason Deployment dialog appears.



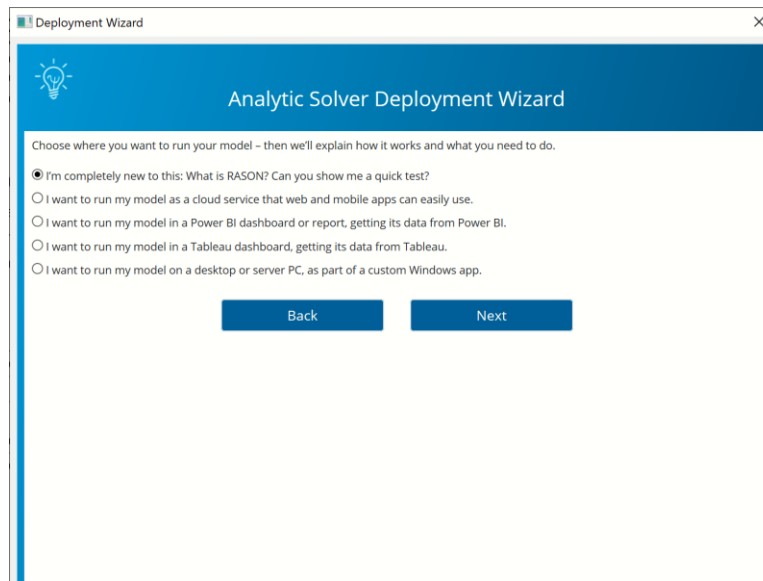
Follow the steps below to convert an optimization, simulation or data science model into a RASON model and solve it in the Cloud, Tableau, Power BI, with Solver SDK or on an automatically created web page.

## Deployment Wizard

If you just aren't sure where to start or what you need, use the handy deployment wizard to help guide you through the RASON Deployment menu options.



When you click "Next" you'll be guided to select the best way to deploy your Excel model. In each case, you or a developer can extend your deployed application by writing C#, JAVA, C++, Python, R, JavaScript, or you can revise your model in Excel and "plug" the new version into the application.



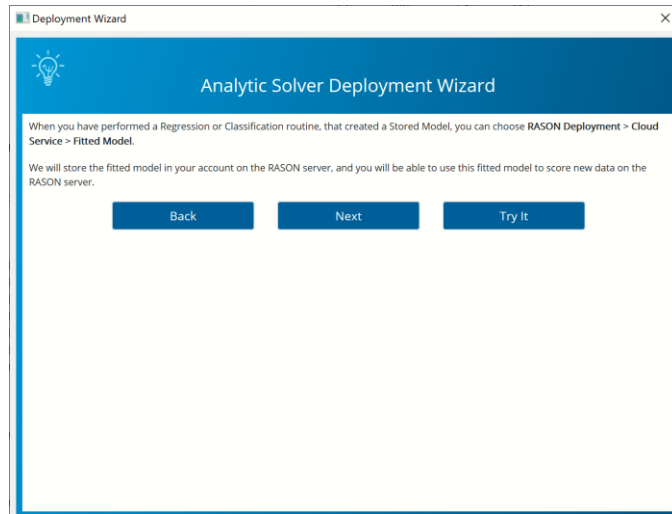
## RASON Deployment Menu

Clicking the second option, "I want to run my model as a cloud service..." gives you four options, RASON Model, Excel Model, Fitted Model and Probability Model.

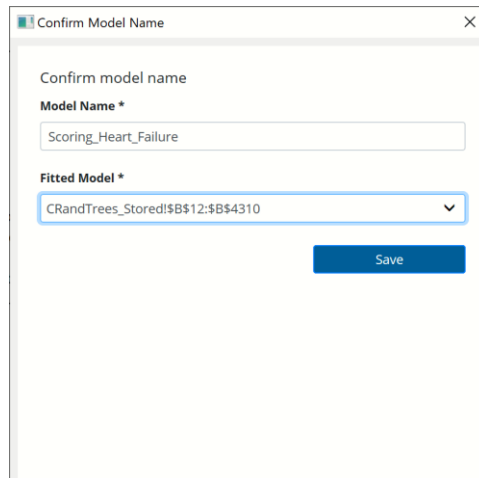
- RASON Model translates your optimization, simulation, data science and decision table model into Rason and POSTs the model to your RASON account.
- Excel Model POSTs the Excel workbook model to the RASON account.
- Fitted Model POSTs the data science or forecasting fitted model saved to an Excel workbook to the RASON Server.
- Probability Model allows the deploying and sharing of probability models, following the open Probability Management 3.0 standard. See the Analytic Solver User Guide for more information.

From here you can maintain and solve your model and get the results on [www.RASON.com](http://www.RASON.com).

Clicking "Next" a few times brings you to the Fitted Model deployment page.



Click "Try It" to deploy your data science/forecasting fitted model to a Web application. If a fitted model is contained in your workbook, the following dialog will appear. Here you can give the RASON model a name and also select the range of the fitted model.



Clicking the third option, "I want to run my model in a Power BI dashboard...", gives you the ability to turn your Excel-based optimization, simulation or data science fitted model into a **Microsoft Power BI Custom Visual**. First, you

simply select rows or columns of data to serve as changeable parameters, next, you select **RASON Deployment – Power BI – Managed Model/Embedded Model**, and thirdly, you save the file created by Analytic Solver. Afterwards, you click the Load Custom Visual icon in Power BI and select the file you just saved. What you get isn't just a chart – it's your *full optimization, simulation or data science fitted model*, ready to accept Power BI data, **run on demand** on the web, and display visual results in Power BI! You simply need to drag and drop appropriate Power BI datasets into the “well” of inputs to match your model parameters. This is possible because Analytic Solver translates your Excel model into **RASON®** then “wraps” a JavaScript-based Custom Visual around the RASON model. For more information on this feature, see the next chapter “Creating Power BI Custom Visuals”.

Clicking the fourth option, "I want to run my model in a Tableau dashboard...", allows you to turn your Excel-based optimization, simulation or data science fitted model into a **Tableau Dashboard Extension**. You simply select rows or columns of data to serve as changeable parameters, then choose **RASON Deployment – Tableau – Managed Model/Embedded Model**, and save the file created by Analytic Solver. In Tableau, you'll see the newly created file under **Extensions** on the left side of the dashboard, where you can drag it onto your dashboard. You'll be prompted to match the parameters your model needs, with data in Tableau. Much like with Power BI, what you get isn't just a chart – it's your *full optimization, simulation or data science fitted model*, ready to accept Tableau data, **run on demand** (using our **RASON** server), and display visual results in Tableau! Note: This feature works (only) with Tableau version 2018.2 or later. For more information on this feature, see the chapter “Creating Custom Extensions in Tableau”.

Continue reading for a step-by-step example that illustrates how to use each of these three RASON Deployment functionality.

## Conversion Exceptions

There are several types of models that Frontline Excel Solvers are not able to convert to RASON models, i.e. Excel models containing errors such as #NUM, #VALUE, etc. For all unsupported features, the conversion generator returns an appropriate message.

Note: When converting a simulation model into the RASON modeling language note that only statistics with explicit numeric arguments such as =PsiPercentile(cell, 0.50) or =PsiPercentile(cell, A1) where A1 = 0.05 are supported. Expressions passed to arguments such as =PsiPercentile(cell, A1+A2), where A1 = .1 and A2 = .3 are not supported.

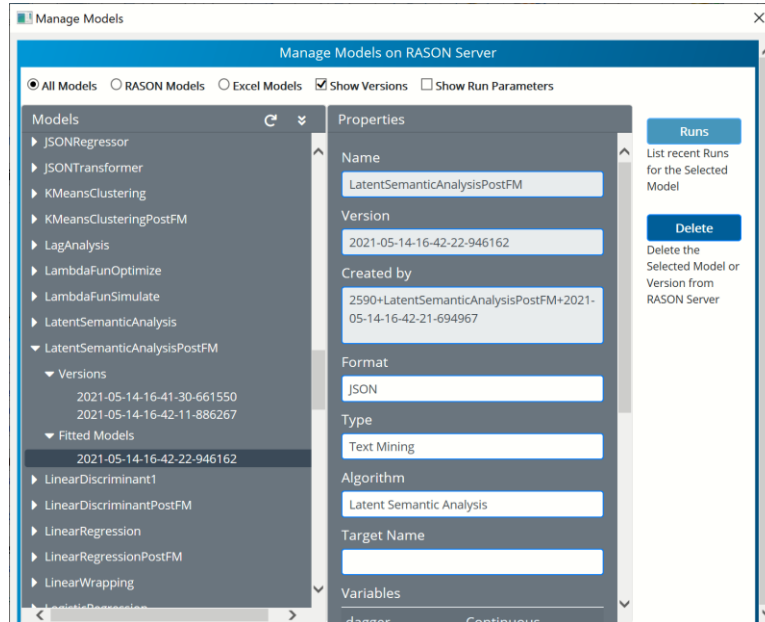
## Manage Models

Click Manage Models (top right) to display your models and model versions from your account on [www.rason.com](http://www.rason.com).

- Selecting a model name in the left column displays the model properties in the right.
- Select all models to see all your models on your Rason account, Rason as well as Excel models.
- Click Show Versions to see all versions of each named model.
- Click Show Run Parameters to display query parameters.

- With a model name selected, clicking the Runs button displays a list of the most recent model instances in the right column.
- Clicking Delete deletes the selected named model (all versions) or a specific version.

The screenshot below displays a fitted model in the Manage Model dialog.



## Cloud Service

Clicking the first option on the RASON Deployment dialog, Cloud Service, gives you four options, RASON Model, Excel Model, Fitted Model and Probability Model.

- *RASON Model* translates the Excel model into Rason and POSTs the model to your RASON account. In Analytic Solver Data Science, this option is applicable only for workflows. Please see the previous chapter for more information.
- *Excel Model* POSTs the Excel workbook model to the RASON account. This option is not applicable for data science/forecasting models.
- *Fitted Model* POSTs the fitted model contained within your Excel workbook to your RASON account.
- *Probability Model* - This option is not supported for data science/forecasting models.

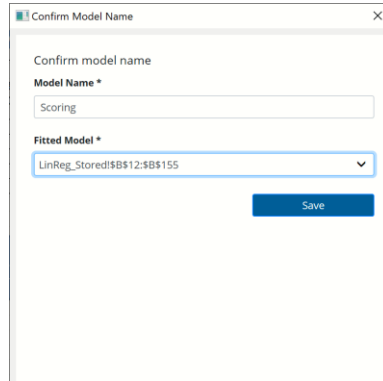
### Cloud Service Fitted Model

The majority of Analytic Solver Data Science routines generate fitted, or stored, models which are used to score (classify, predict or forecast) new data points. Selecting Cloud Service – Fitted Model POSTs the fitted model to your account on the RASON server allowing users to utilize their fitted models to score new

data on the RASON Server. Continue reading to follow a step-by-step example illustrating this feature.

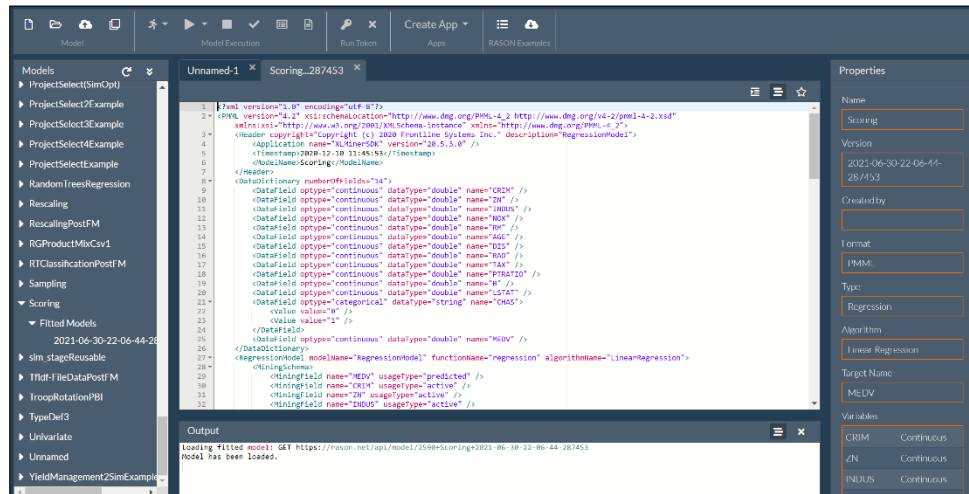
On the Analytic Solver ribbon, click Help – Example Models – Forecasting / Data Science Examples. Scroll down and click the Scoring example link to open Scoring.xlsx. This workbook contains a stored data science model, LinReg\_Stored.

Click RASON Deployment – Cloud Service – Fitted Model. The Confirm Model Name dialog opens.

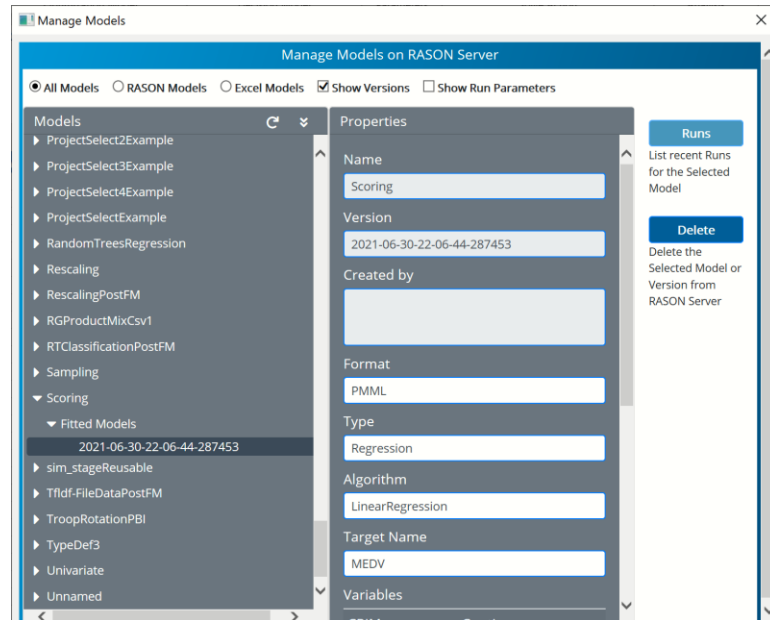


- Enter a Model name. This example will use the default, Scoring, but any name will work.
- Click the down arrow beneath Fitted Model to select the fitted model contained within the workbook, LinReg\_Stored!\$B\$12:\$B\$155.
- Click Save to POST the fitted model to your RASON account.

From here you can maintain the fitted model through RASON Decision Services on RASON.com.



The fitted model will also be listed in Manage Models.



## Power BI

Clicking the second option on the RASON Deployment dialog, Power BI, gives you three options, Managed Model, Embedded Model and Fitted Model. *Fitted Model* POSTs the data science or forecasting fitted model contained within your Excel workbook to your RASON account. The first two options, Managed Model and Embedded Model, do not support data science/forecasting models.

### Power BI Fitted Model

Users can choose to deploy and share data science and machine learning models, trained in Analytic Solver or RASON, to the Azure cloud, and use them directly for classification and prediction (without needing auxiliary “code” in R or Python, RASON or Excel). (In the past, Analytic Solver Data Science models could only be scored from within Excel, either online or desktop.) This example illustrates how one could use a fitted model created in Analytic Solver Data Science to score data residing in Power BI, Microsoft’s data analytics platform.

Open the Scoring\_Heart\_Failure.xlsx example workbook by clicking Help – Example Models – Forecasting / Data Science Examples on the Analytic Solver ribbon. This example uses the [Heart Failure Clinical Records Dataset](#)<sup>2</sup>, which contains thirteen variables describing 299 patients experiencing heart failure. The [journal article](#) referenced here discusses how the authors analyzed the dataset to first rank the features (variables) by significance and then used the Random Trees machine learning algorithm to fit a model to the dataset to be used to classify patients at risk of death due to heart failure.

<sup>2</sup> Davide Chicco, Giuseppe Jurman: Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone. BMC Medical Informatics and Decision Making 20, 16 (2020). ([link](#))



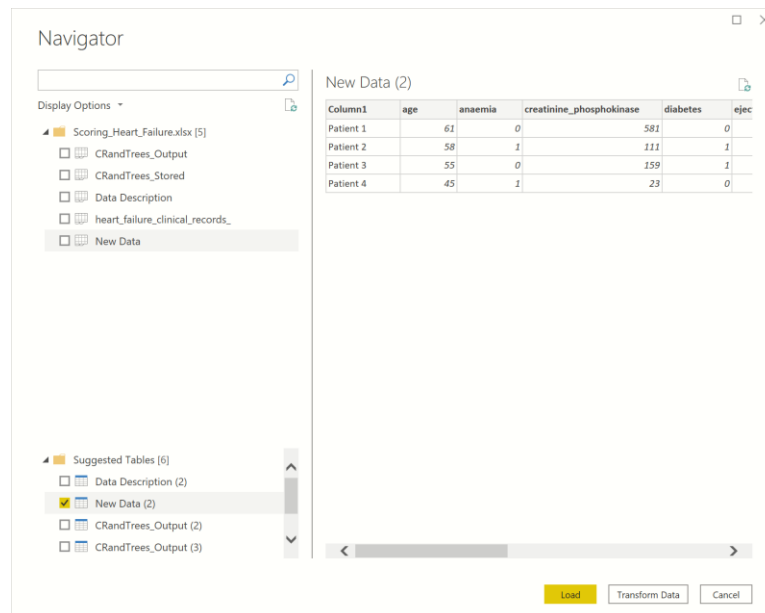
DEATH\_EVENT is the output variable in the Random Trees Classification model. This variable can either be a 1 or 0. A patient classified as a 1 is predicted to perish due to complications of heart failure.

See the Feature Selection chapter within the Data Science User Guide for an example that attempts to emulate the results of that study. For now, this example dataset will illustrate how to score new data contained within Power BI using the Random Trees Classification fitted model saved in the workbook on the CRandTrees\_Stored worksheet.

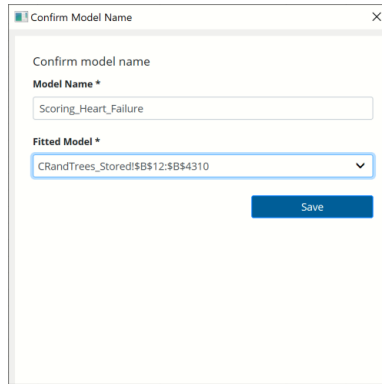
Click the New Data tab to open a list of new patient records which require classification using the stored model, CRandTrees\_Stored. It is imperative that patients with a classification of "1" for the output variable, DEATH\_EVENT, receive medical mitigation in order to reduce the chance of death from heart failure.

Open Power BI and import the New Data table. (Open Power BI and click Get Data – Excel and select the New Data table, then click Load.)

Note: This dataset includes four rows of new data but only 1 row is required for scoring.

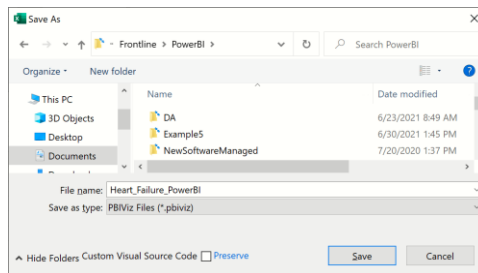


Click back to the CRandTrees\_Stored worksheet in Excel and then click RASON Deployment – Power BI – Fitted model to POST the fitted model, CRandTrees\_Stored, to the RASON Server. This fitted model will be used to score the data that was just imported into Power BI.



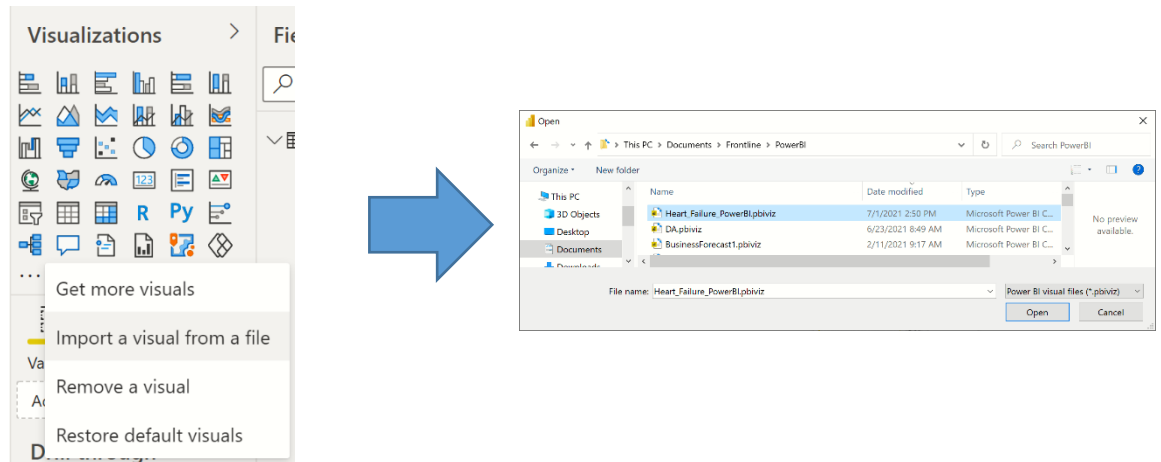
Click the down arrow beneath Fitted Model and select the range CRandTrees\_Stored!\$B\$12:\$B\$4310. This is the range where the fitted model is saved on the CRandTrees worksheet.

Enter a file name for the Power BI custom visual file that will be created and saved. This example uses Heart\_Failure\_PowerBI



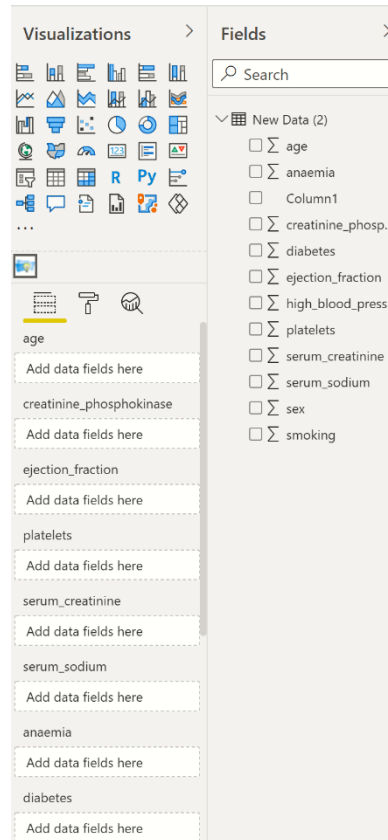
Click Save to save the Power BI Custom Visual. A command prompt window will appear on your screen and, after several moments, disappear. It is during this time that your Excel model is being translated into the **RASON®** modeling language, and your Power BI custom visual is created.

Click back to Power BI and import the Power BI custom visual by clicking the three dots beneath Visualizations, then browse to open Heart\_Failure\_PowerBI.pbviz.

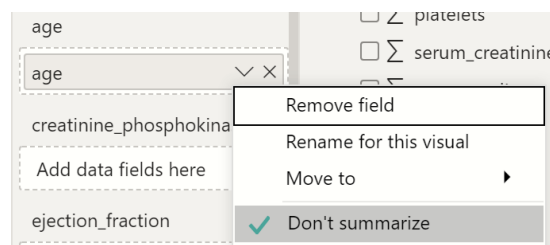


Click the newly created Heart\_Failure\_PowerBI custom visual to display the custom visual data wells.

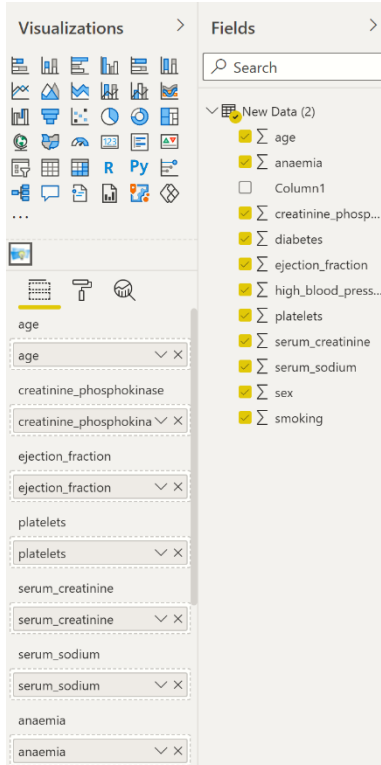
Expand New Data on the right to display the fields contained within the data table.



Match the data into the appropriate data well by dragging the New Data columns to the custom visual data wells, i.e. drag "age" to "age", "anaemia" to "anaemia", etc. Then click the down arrow next to each and select "Don't summarize".



Power BI should look similar to the screenshot below after all fields are matched with the appropriate data wells.



Once the last data column is match to the last data well, the custom visual will score the data and display the results in the Power BI dashboard, as shown in the screenshot below.

Output results contain the data to be scored plus the column containing the classification of the output variable, DEATH\_EVENT.

| DEATH_EVENT scoring results |           |                  |              |         |          |                     |     |         |             |
|-----------------------------|-----------|------------------|--------------|---------|----------|---------------------|-----|---------|-------------|
| ejection fraction           | platelets | serum_creatinine | serum_sodium | anaemia | diabetes | high blood pressure | sex | smoking | DEATH_EVENT |
| 60                          | 166000    | 5.7              | 132          | 0       | 1        | 1                   | 0   | 1       | 1           |
| 30                          | 162000    | 1.7              | 140          | 1       | 1        | 1                   | 1   | 1       | 1           |
| 50                          | 300000    | 9.2              | 116          | 1       | 0        | 0                   | 0   | 0       | 1           |
| 38                          | 263358    | 1.4              | 137          | 0       | 0        | 0                   | 1   | 0       | 0           |

The scoring results indicate that the first three, out of the four patients, require some sort of medical mitigation in order to decrease the risk of death from heart failure.

Note that the order that Power BI scores the data may be slightly different than the order that Analytic Solver Data Science scores the data.

## Tableau

Clicking the third option on the RASON Deployment dialog, Tableau, gives you three options, Managed Model, Embedded Model and Fitted Model. *Fitted Model* POSTs the data science or forecasting fitted model contained within your Excel workbook to your RASON account. The first two options, Managed Model and Embedded Model, do not support data science models.

## Tableau Fitted Model

As introduced in the previous section, What's New in Analytic Solver V2021.5, users can now deploy and share data science and machine learning models, trained in Analytic Solver or RASON, to the Azure cloud, and use them directly for classification and prediction (without needing auxiliary “code” in R or Python, RASON or Excel). (In the past, Analytic Solver Data Science models could only be scored from within Excel, either online or desktop.) This example illustrates how one could use a fitted model created in Analytic Solver Data Science to score data residing in Tableau, a visualization platform. For more information on Tableau, see the chapter Creating Custom Extensions in Tableau that appears later in this guide.

Open the Scoring\_Heart\_Failure.xlsx example workbook by clicking Help – Example Models – Forecasting / Data Science Examples on the Analytic Solver ribbon. This example uses the [Heart Failure Clinical Records Dataset](#)<sup>3</sup>, which contains thirteen variables describing 299 patients experiencing heart failure. The [journal article](#) referenced here discusses how the authors analyzed the dataset to first rank the features (variables) by significance and then used the Random Trees machine learning algorithm to fit a model to the dataset to be used to classify patients at risk of death due to heart failure.

DEATH\_EVENT is the output variable in the Random Trees Classification model. This variable can either be a 1 or 0. A patient classified as a 1 is predicted to perish due to complications of heart failure.

See the Feature Selection chapter within the Analytic Solver Data Science User Guide for an example that attempts to emulate the results of that study. For now, this example dataset will illustrate how to score new data contained within Tableau, using the Random Trees Classification fitted model saved in the workbook.

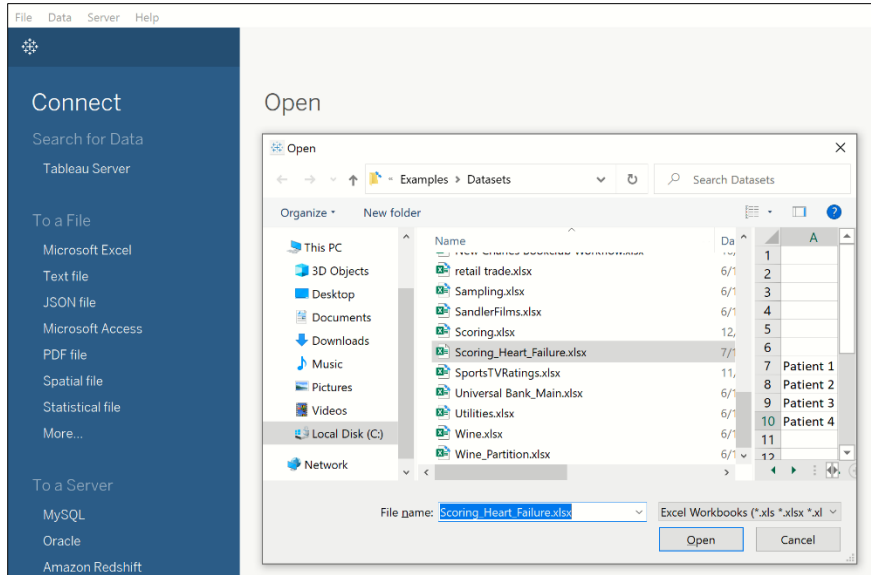
Click the New Data tab to open a list of new patient records which require classification using the stored model, CRandTrees\_Stored. This classification is required in order to determine which patients require medical mitigation to reduce the chance of death from heart failure.

Open Tableau and import the New Data table.

Open Tableau and click Connect To a File: Microsoft Excel. Browse to the location of the Scoring\_Heart\_Failure.xlsx file typically C:\ProgramData\Frontline Systems\Datasets, then click Open.

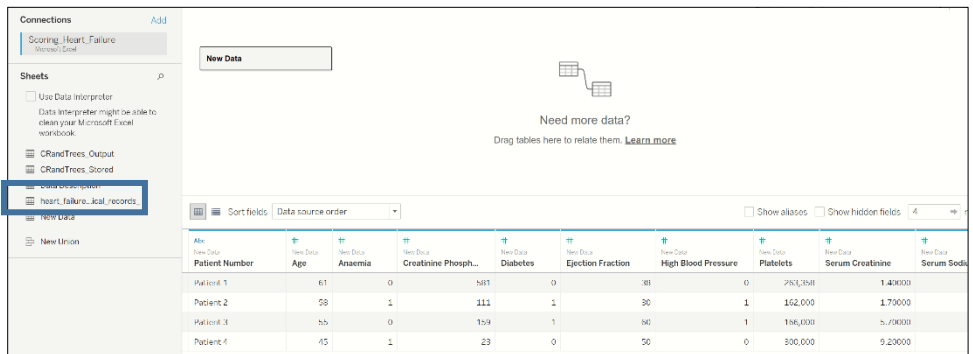
---

<sup>3</sup> Davide Chicco, Giuseppe Jurman: Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone. BMC Medical Informatics and Decision Making 20, 16 (2020). ([link](#))

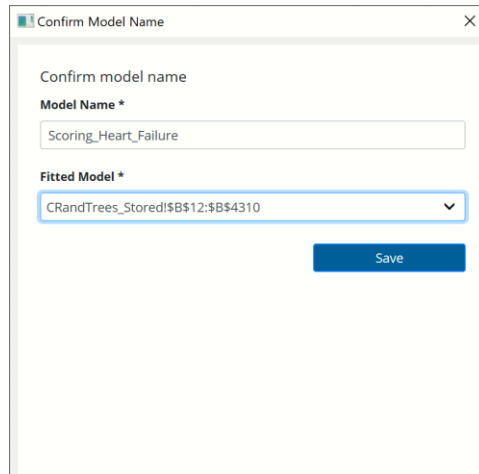


Click the New Data, on the bottom left, to open the table in Tableau. This table contains a list of new patient records which require classification using the stored model, CRandTrees\_Stored. It is imperative that patients with a classification of "1" for the output variable, DEATH\_EVENT, receive medical mitigation in order to reduce the chance of death from heart failure.

Note: This dataset includes four rows of new data but only 1 row is required for scoring.

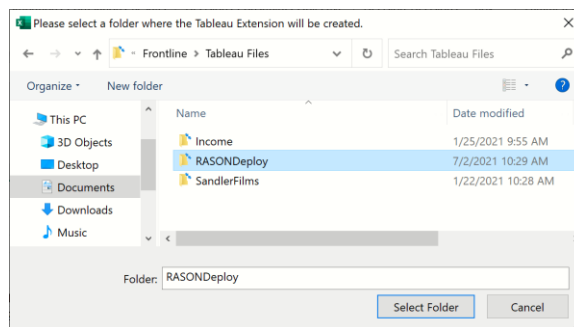


Click back to the CRandTrees\_Stored worksheet in Excel and then click RASON Deployment – Tableau– Fitted model to POST the fitted model, CRandTrees\_Stored, to the RASON Server. This fitted model will be used to score the data that was just imported into Tableau

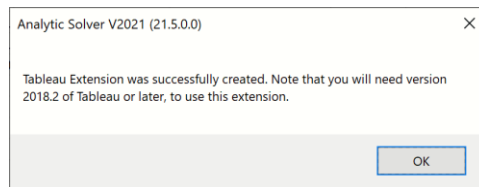


Click the down arrow beneath Fitted Model and select the range CRandTrees\_Stored!\$B\$12:\$B\$4310. This is the range where the fitted model is saved on the CRandTrees worksheet.

Select a folder to save the Tableau custom extension. This example uses the RASONDeploy folder.



Click Select Folder. It is during this time that your Excel model is being translated into the **RASON®** modeling language, and your Tableau custom extension is created. After a few moments, a dialog will appear indicating the Tableau extension has been created successfully.



### ***Starting up the Server***

Since Tableau extensions are simply web pages, we will first need to start up a web server to serve our content. For this example, we will serve up the webpage to the default location. To do so, open a command prompt, navigate to the root of the extensions repository and run “http-server -p 8000”.

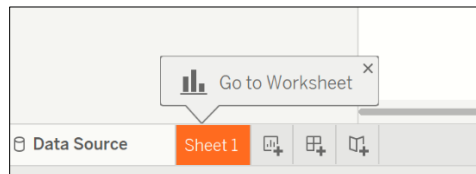
```

http-server
Starting up http-server, serving ./
Available on:
 http://192.168.86.151:8000
 http://127.0.0.1:8000
Hit CTRL-C to stop the server
[2021-07-02T18:11:31.091Z] "GET /" "Tableau Desktop 20212.21.0605.1023; pro; libcurl-client; 64-bit; en_US; Microsoft Windows 10 Pro (Build 19042); Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) QtWebEngine/5.15.0 Chrome/80.0.3987.163 Safari/537.36"
(node:17576) [DEP0066] DeprecationWarning: OutgoingMessage.prototype._headers is deprecated
[2021-07-02T18:11:31.103Z] "GET /favicon.ico" "Tableau Desktop 20212.21.0605.1023; pro; libcurl-client; 64-bit; en_US; Microsoft Windows 10 Pro (Build 19042); Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) QtWebEngine/5.15.0 Chrome/80.0.3987.163 Safari/537.36"
[2021-07-02T18:11:31.111Z] "GET /favicon.ico" Error (404): "Not found"
[2021-07-02T18:11:31.135Z] "GET /tableau.js" "Tableau Desktop 20212.21.0605.1023; pro; libcurl-client; 64-bit; en_US; Microsoft Windows 10 Pro (Build 19042); Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) QtWebEngine/5.15.0 Chrome/80.0.3987.163 Safari/537.36"
[2021-07-02T18:11:31.138Z] "GET /solver.js" "Tableau Desktop 20212.21.0605.1023; pro; libcurl-client; 64-bit; en_US; Microsoft Windows 10 Pro (Build 19042); Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) QtWebEngine/5.15.0 Chrome/80.0.3987.163 Safari/537.36"
[2021-07-02T18:11:31.147Z] "GET /solver.css" "Tableau Desktop 20212.21.0605.1023; pro; libcurl-client; 64-bit; en_US; Microsoft Windows 10 Pro (Build 19042); Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) QtWebEngine/5.15.0 Chrome/80.0.3987.163 Safari/537.36"
>

```

This command starts up a simple http server listening on port 8000.

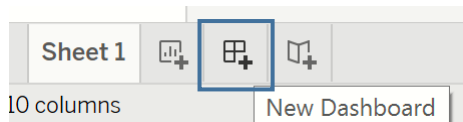
Click back to Tableau and click Sheet 1.



Drag Patient Number to Rows and then drag the 11 features (or variables) to Columns.

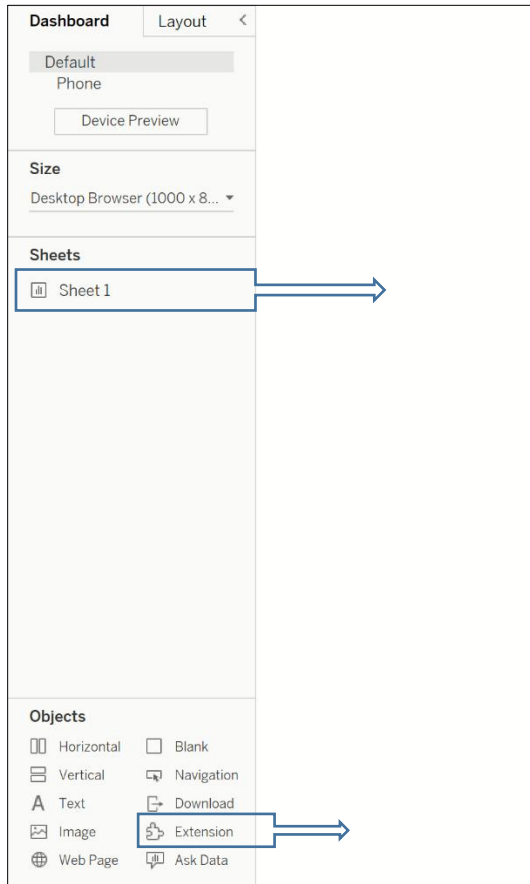


Then click New Dashboard to open a new Tableau dashboard.

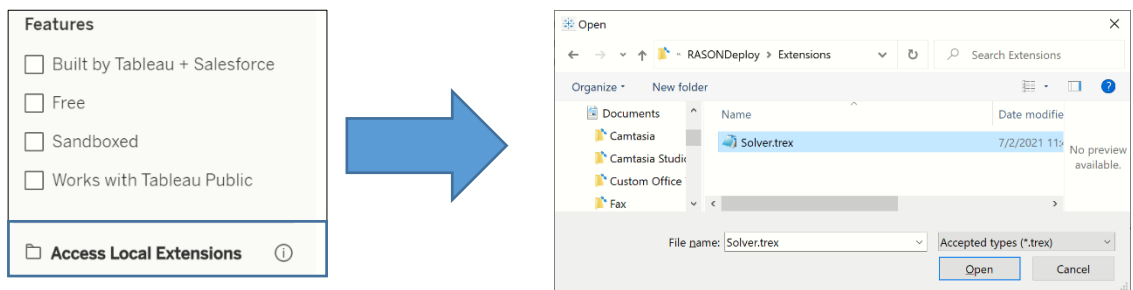




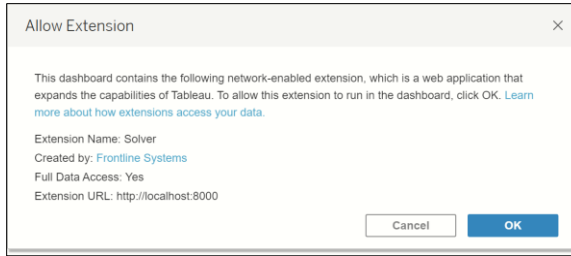
Once opened, drag Sheet 1, under Sheets, and Extension, under Objects to the dashboard.



The "Add an Extension" dialog opens. On this dialog, click **Access Local Extensions** on the bottom left and browse to the location of the Tableau .trex file.

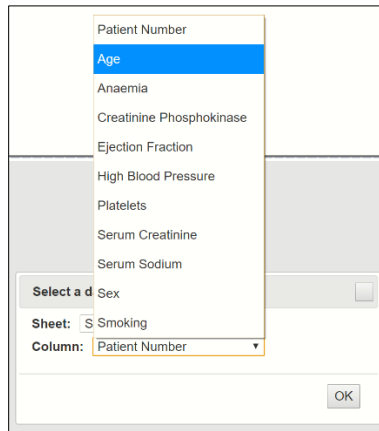


Click Open to open the Tableau extension. Then click OK to allow the extension to run the Tableau dashboard.

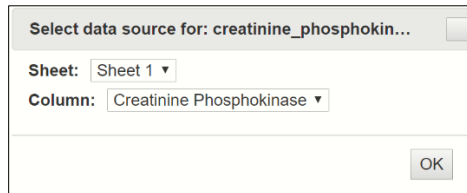


A series of dialogs will appear which allow the dimensions from the Tableau worksheet to be matched with the dimensions in the Solver custom Tableau extension (Solver.trex).

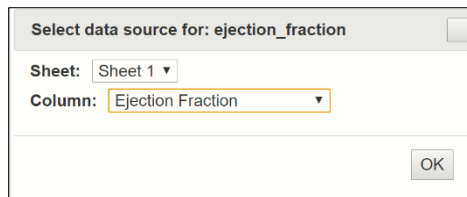
- Match the data source "age" dimension with the worksheet "age" dimension.



- Match the data source "creatinine\_phosphokinase" dimension with the worksheet "Creatinine Phosphokinase" dimension.



- Match the data source "ejection\_fraction" dimension with the worksheet "Ejection Fraction" dimension.



- Match the data source "Platelets" dimension with the worksheet "Platelets" dimension.

Select data source for: platelets

Sheet: Sheet 1

Column: Platelets

OK

- Match the data source "serum\_creatinine" dimension with the worksheet "Serum Creatinine" dimension.

Select data source for: serum\_creatinine

Sheet: Sheet 1

Column: Serum Creatinine

OK

- Match the data source "serum\_sodium" dimension with the worksheet "Serum Sodium" dimension.

Select data source for: serum\_sodium

Sheet: Sheet 1

Column: Serum Sodium

OK

- Match the data source "anaemia" dimension with the worksheet "Anaemia" dimension.

Select data source for: anaemia

Sheet: Sheet 1

Column: Anaemia

OK

- Match the data source "diabetes" dimension with the worksheet "Diabetes" dimension.

Select data source for: diabetes

Sheet: Sheet 1

Column: Diabetes

OK

- Match the data source "high\_blood\_pressure" dimension with the worksheet "High Blood Pressure" dimension.

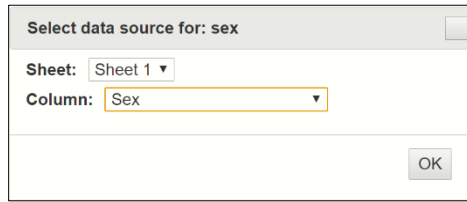
Select data source for: high\_blood\_pressure

Sheet: Sheet 1

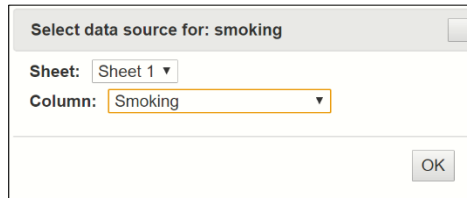
Column: High Blood Pressure

OK

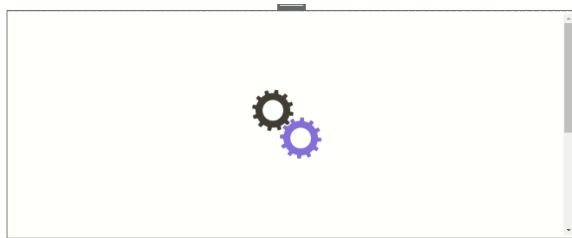
- Match the data source "sex" dimension with the worksheet "Sex" dimension.



- Match the data source "smoking" dimension with the worksheet "Smoking" dimension.



Once the last dimension is matched, Tableau immediately begins generating the custom extension.



The screenshot below displays the completed custom extension.

Output results contain the data to be scored plus the column containing the classification of the output variable, DEATH\_EVENT.

| DEATH_EVENT_scoring results |                  |              |         |          |                     |     |         |             |
|-----------------------------|------------------|--------------|---------|----------|---------------------|-----|---------|-------------|
| s                           | serum_creatinine | serum_sodium | anaemia | diabetes | high_blood_pressure | sex | smoking | DEATH_EVENT |
|                             | 1.4              | 137          | 0       | 0        | 0                   | 1   | 0       | 0           |
|                             | 1.7              | 140          | 1       | 1        | 1                   | 1   | 1       | 1           |
|                             | 5.7              | 132          | 0       | 1        | 1                   | 0   | 1       | 1           |
|                             | 9.2              | 116          | 1       | 0        | 0                   | 0   | 0       | 1           |

The scoring results indicate that the first three, out of the four patients, require some sort of medical mitigation in order to decrease the risk of death from heart failure.

# Bringing Big Data into Excel Using Apache Spark

---

## Introduction

Large amounts of data are being generated and collected continuously from a multitude of sources every minute of every day. From your toothbrush to your vehicle GPS to Twitter/Facebook/Google/Yahoo, data is everywhere. Being able to make decisions based on this information requires the ability to extract trends and patterns<sup>2</sup> that can be buried deeply within the numbers.

Generally these large datasets contain millions of records (rows) requiring multiple gigabytes or terabytes of storage space across multiple hard drives in an external compute cluster. Analytic Solver Data Science enables users to ‘pull’ sampled and summarized data into Excel from compute clusters running Apache Spark, the open-source software widely embraced by Big Data vendors and users.

---

## Sampling and Summarizing Big Data

This example illustrates how to use the Big Data Sample/Summarization feature using Data stored across an Apache Spark compute cluster where the Frontline Systems access server is installed. By drawing a *representative sample* of Big Data from all the nodes in the cluster, Excel users can easily train data science and text mining models directly on their desktops.

In this example, we will use the Airline dataset. The data used in this example consists of flight arrival and departure information for all commercial flights within the USA dating from October 1987 to April 2008. This data was obtained from 29 commercial airlines and 3,376 airports and consists of 3.2 million cancelled flights and 25 million flights at least 15 minutes late. This is a large dataset with nearly 120 million records requiring 1.6 GB of storage space when compressed and 12 GB of storage space when uncompressed. Data was obtained from the Research and Innovative Technology Administration (RITA) which coordinates the U.S. Department of Transportation research programs. Note: Southwest (WN), American Airlines (AA), United Airlines (UA), US Airways (US), Continental Airlines (CO), Delta Airlines (DL), Northwest Airlines (NW) and Alaska Airlines (AS) are the only airlines where data is available for all 20 years. Recall the annual revenue from the domestic airline industry is \$157 billion. This public dataset was obtained from [here](#). Navigate to this webpage to explore details about this dataset. For supplemental data including the location of each airport, plane type and meteorological data pertaining to each flight, click [here](#).

The information contained in this large dataset could allow us to answer or understand the following questions or issues:

- What are the airports most prone to departure delays? What airports tend to have the most arrival delays?

- What are the times of day and days of week that are most susceptible to departure/arrival delay?
- How can we understand flight patterns as they respond to well-known events? (i.e., examining the data before and after September 2011)
- How many miles per year does each plane by carrier fly?
- When is the best time of day/day of week/time of year to fly to minimize delays?
- How does the number of people flying between different locations change over time?
- How well does weather predict plane delays?
- Can you detect cascading failures as delays in one airport creates delays in others? Are there critical links in the system?
- Understanding flight patterns between the pair of cities that you fly between most often, or all flights to and from a major airport like Chicago (ORD)
- Average arrival delay in minutes by flight or by year?
- How many flights were cancelled, at least 15 minutes late, etc?
- How many flights were less than 50 miles?

## Connecting to an Apache Spark Cluster

The first step in connecting Analytic Solver Data Science to your organization's own Apache Spark cluster is to contact Frontline Systems Sales and Technical Support at 775-831-0300. After the server-side software package is installed, the proper entries for the cluster options can be entered as shown in the example below.

## Storage Sources and Data Formats

Analytic Solver Data Science can process data from [Hadoop Distributed File System \(HDFS\)](#), local file systems that are visible to Spark cluster, and [Amazon S3](#). Performance is best with HDFS, and it is recommended that you load data from a local file system or Amazon S3 into HDFS. If the local file system is used, the data must be accessible at the same path on all Spark workers, either via a network path, or because it was copied to the same location on all workers.

At present, Analytic Solver Data Science can process data in [Apache Parquet](#) and CSV (delimited text) formats. Performance is far better with Parquet, which stores data in a compressed, columnar representation; it is highly recommended that you convert CSV data to Parquet before you seek to sample or summarize the data.

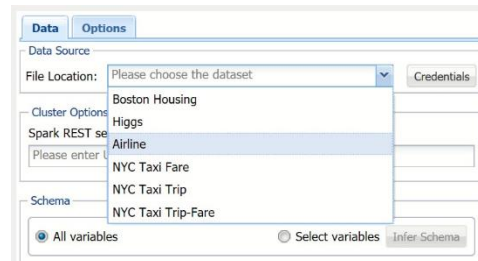
## Sampling from a Large Dataset

If using desktop Analytic Solver Data Science, click **Get Data – Big Data – Sample** to open the *Sample Big Data* dialog. Enter the location of the file for *File Location* and the URL for the Spark Server for *Spark REST server URL*. This example uses the Airline dataset (described above) installed on a Frontline

operated Apache Spark cluster. If your dataset is located on Amazon S3, click Credentials to enter your Access and Secret Keys.

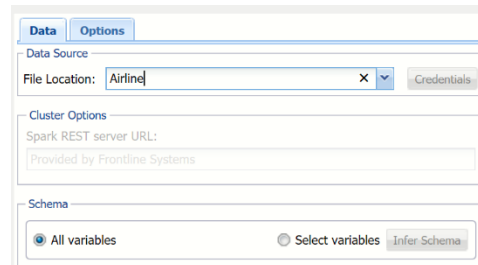
If using AnalyticSolver.com, click the down arrow to open the *File Location* drop down menu and select **Airline**.

**Note: Frontline Systems no longer offers Big Data datasets, including the Airline dataset, as a free service. However, the instructions below are applicable to any connected Spark cluster.**



Keep *All variables* selected to include all variables in the dataset in the sample.

To only include specific variables, you would choose *Select variables*, then click *Infer Schema*. All variables contained in the dataset will be reported under *Variables*. Use the >/< and >>/<< buttons to select variables to be included in the sample.



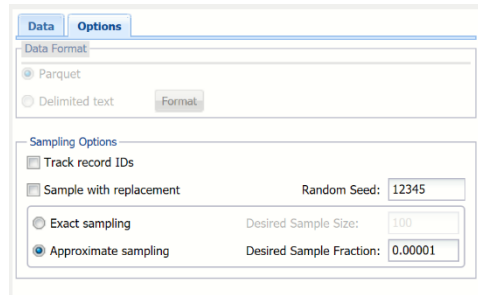
Click the **Options** tab.

Select **Approximate Sampling**. When this option is selected, the size of the resultant sample will be determined by the value entered for *Desired Sample Fraction*. Approximate sampling is much faster than Exact Sampling. Usually, the resultant fraction is very close to the Desired Sample Fraction so this option should be preferred over exact sampling as often as possible. Even if the resultant sample slightly deviates from the desired size, this would be easy to correct in Excel.

Enter **0.00001** for *Desired Sample Fraction*. This option controls the fraction of the total number of records in the full dataset that is expected to be included in the generated sample. Since our dataset contains about 120 million records, our sample will contain approximately 1,200 records. If *Sampling with Replacement* is selected, the value for *Desired Sample Fraction* is the expected number of times each record can be chosen and must be greater than 0. If Sampling **without** replacement (i.e. Sampling with Replacement is not selected), the Desired Sample Fraction becomes the probability that each element is chosen and, as a result, *Desired Sample Fraction* must be between 0 and 1.

Keep *Random Seed* at the default of **12345**. This value initializes the random number generator. This option allows you to generate reproducible samples. *Track record IDs* and *Sample with replacement* should remain unchecked.

Please see the Analytic Solver Data Science Reference guide for a complete description of each option included on this dialog.



Clicking *Submit* sends a request for sampling to the compute cluster but does not wait for completion. The result is output containing the Job ID and basic information about the submitted job so multiple submissions may be identified. This information can be used at any time later for querying the status of the job and generating reports based on the results of the completed job.

Clicking *Run* sends a request for sampling to the compute cluster and waits for the results. Once the job is completed and results are returned to the Analytic Solver client, a report is inserted into the Model tab on the Data Science task pane under Reports – Text Science containing the sampling results.

Click **Submit**. Results will be inserted into the Data Science Task Pane under Results – Sampling – Run 1. Open BD\_Sampling.

**Output Navigator**

Inputs    Submission Results

**Elapsed Times in Milliseconds**

| Data Reading Time | Algorithm Time | Report Time | Total |
|-------------------|----------------|-------------|-------|
| 0                 | 0              | 314         | 314   |

**Inputs**

| Data Info          |                                         |
|--------------------|-----------------------------------------|
| Data location type | Cluster: Hadoop distributed file system |
| Data format        | Apache Parquet                          |
| Dataset            | Airline                                 |
| Credentials        | Unrestricted                            |

| Cluster Configuration          |                               |
|--------------------------------|-------------------------------|
| Apache Spark REST server URL   | Provided by Frontline Systems |
| Total number of CPU cores used | 0                             |

**Big Data Sampling Parameters**

|                         |             |
|-------------------------|-------------|
| Track record IDs?       | FALSE       |
| Sampling type           | Approximate |
| Desired sample fraction | 0.00001     |
| With replacement?       | FALSE       |
| Random seed             | 12345       |

**Submission Results**

|        |                                      |
|--------|--------------------------------------|
| Status | Started successfully                 |
| Job ID | e4df64c9-1d4c-40a0-a41b-e7720c179372 |

This report displays the details about the chosen dataset, selected options for sampling and the job identifier required for identifying the submission on the cluster.

Click **Get Data – Big Data – Get Results** on the Data Science ribbon to open the *Big Data: Get Results* dialog. Click the down arrow to the right of *Job identifier* and select the previously submitted job. Click **Get Info** to obtain the status of the Job from the cluster.





The *Inputs* section displays the information about the dataset, cluster configuration, details on the running time, the options chosen during setup, and a summary of the data including the size and dimensionality of the full and sampled datasets. Since *Approximate Sample* was selected, we can expect the resulting fraction to be slightly different from the desired fraction. In our example, the resulting fraction, approximately 0.00001014 (1,184/116,701,402) is very close to the requested fraction (0.00001). (Recall *Approximate Sample* was selected and *Desired Sample Fraction* was entered on the *Options* tab during setup.)

Scroll down to see the full and sampled data schemas. Since we choose to include all variables in the sample, the set of columns in the full and sampled datasets is the same.

Further down we see the *Sampled Data*, which includes 1,184 records, as indicated by the *Number of records – sample* field under *Sampled Data Summary*.

Now that the representative data sample has been drawn and is available in the output, all of the methods and features included in *Analytic Solver Data Science* are literally at your fingertips. We could choose to explore our sampled data by creating visualizations using the *Chart Wizard*, transform the data using *Analytic Solver Data Science's* data transformation utilities, build classification/prediction models to forecast arrival and departure times, predict airport delays, estimate total flight times and perform any other analytic tasks that can address numerous challenges that *Big Data*, and the *Airline* dataset in particular, present to the data scientists and analysts.

## Summarizing a Large Dataset

The *Big Data Summarization* feature in *Analytic Solver Data Science* translates the lightning-fast cluster computing capabilities from the state of the art *Big Data* engine, *Apache Spark*, to the simple and easy to use, "point & click" interface within *Excel*. This very powerful yet intuitive tool is useful for rapid extraction of key metrics contained in data, which can be immediately used by data analysts and decision makers. The new *Summarization* feature in *Analytic Solver Data Science* provides similar functionality as standard *SQL* engines, but for the data, volume and complexity which extends far beyond your desktop or laptop computer. This tool is a great assistant for composing reports, constructing informative visualizations, building prescriptive and predictive models that can drive the directions of consequent analysis.

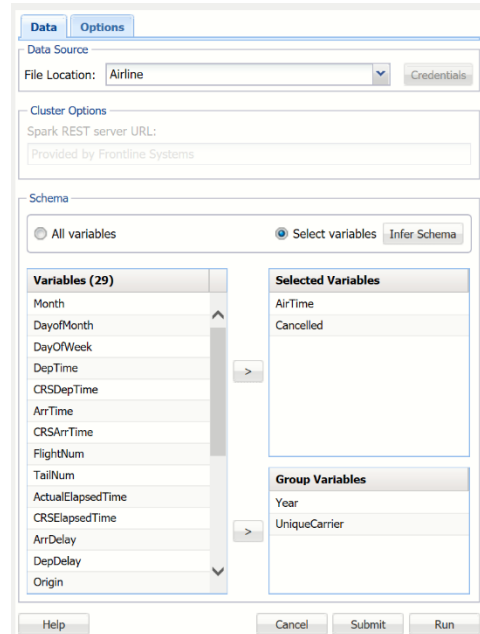
Now we will illustrate how to utilize this easy-to-use yet powerful tool by summarizing the *Airport* dataset and using the information obtained to answer the following three questions:

- What carrier has the most domestic flights by year?
- Who are the most reliable airlines?
- Who are the least reliable airlines?

Click **Get Data – Big Data – Summarize** on the *Data Science Ribbon*. This time we will select a subset of variables for summarization along with grouping variables, so (after entering the *File location*, *Spark REST server URL* and *file Credentials*, if needed) choose **Select variables** and click **Infer Schema**. Afterwards transfer **ArrTime**, and **Cancelled** to the *Selected Variables* grid and **Year** and **UniqueCarrier** to the *Group Variables* grid. *Group Variables* are variables from the dataset that are treated as key variables for aggregation. In

this example, the variables will be grouped so that all records with the same Year and UniqueCarrier are included in the same group, and then all aggregate functions for each group will be calculated.

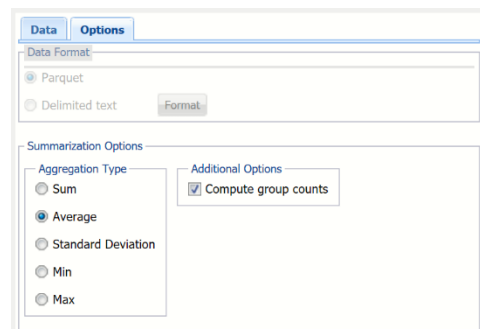
Note: If *All variables* is selected, the result is a simple aggregation of all variables across the entire dataset which can be used to quickly obtain overall statistics.



Click the **Options** tab and select **Average** for *Aggregation Type* and **Compute Group Counts**.

*Aggregation Type* provides 5 statistics that can be inferred from the dataset: sum, average, standard deviation, minimum and maximum.

The option *Compute group counts* is enabled when 1 or more Grouping Variables is selected. When this option is selected, the number of records belonging to each group is computed and reported.

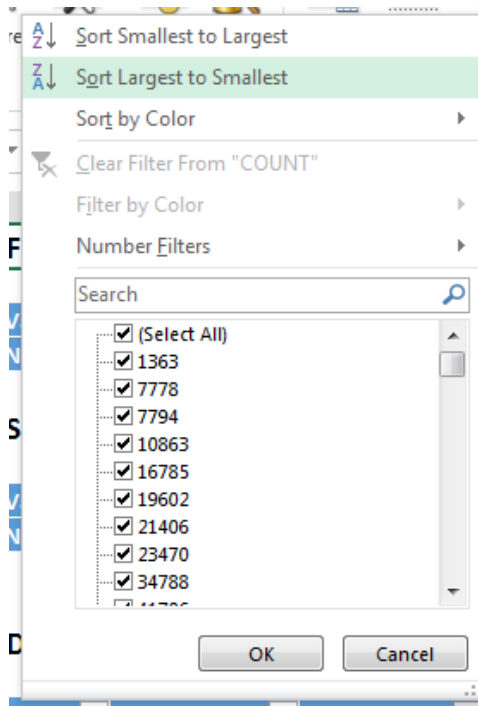


Click **Run** to send a request for a summarization job to the cluster and wait for the results. Once the job is completed, *BD\_Summarization* will be inserted into the current workbook and the Data Science task pane under *Results – Sampling – Run 2*.

Once again the *Inputs* section recaps the dataset details, the cluster configuration, the time taken to complete the job, the options selected during setup and the number of columns and records in both the full dataset and the summarized data. *Full Data Schema* displays all variables in the dataset while *Summarized Data Schema* displays only the variables that were selected during setup.

Scroll down to *Group Counts* to examine the number of records belonging to each Year and UniqueCarrier. In this example, there were 405,598 US Airways flights in 2003 and 684,961 Southwest flights in 1995.

Click the down arrow next to **Count** and select **Sort Largest to Smallest** to answer our first question, "What carrier has the most domestic flights by year?"



Once, the table is sorted on the **Count** column, we see that Southwest (WN) holds the largest market share in domestic flights for years 2005 - 2008.

|    | B                   | C    | D             | E       | F |
|----|---------------------|------|---------------|---------|---|
| 38 | <b>Group Counts</b> |      |               |         |   |
| 39 |                     |      |               |         |   |
| 40 | Group ID            | Year | UniqueCarrier | Count   |   |
| 41 | Group 55            | 2008 | WN            | 1192111 |   |
| 42 | Group 9             | 2007 | WN            | 1158878 |   |
| 43 | Group 248           | 2006 | WN            | 1090370 |   |
| 44 | Group 192           | 2005 | WN            | 1027275 |   |
| 45 | Group 232           | 1990 | US            | 1002485 |   |
| 46 | Group 136           | 2004 | WN            | 980301  |   |
| 47 | Group 80            | 2003 | WN            | 948848  |   |
| 48 | Group 116           | 1997 | DL            | 921831  |   |
| 49 | Group 139           | 1992 | DL            | 916593  |   |
| 50 | Group 171           | 1998 | DL            | 914955  |   |
| 51 | Group 229           | 1999 | DL            | 914080  |   |
| 52 | Group 226           | 2000 | DL            | 907939  |   |
| 53 | Group 283           | 1991 | US            | 907184  |   |
| 54 | Group 213           | 2000 | WN            | 902660  |   |
| 55 | Group 193           | 1993 | DL            | 888896  |   |

Scroll down to Summary Data to find some evidence (that can be further verified) of the most and least "reliable" airlines. Click the down arrow next to *Cancelled\_AVG* and sort from largest to smallest. The airline with the largest average percentage of cancelled flights is Eastern Airlines (EA) with a little over 10% of their flights cancelled on average in 1989. ExpressJet Airlines (EV) and America West (HP) round out the top three spots with 4.5% and 4.3% respectively.

|     | B                   | C         | D    | E             | F           | G             |
|-----|---------------------|-----------|------|---------------|-------------|---------------|
| 329 | <b>Summary Data</b> |           |      |               |             |               |
| 330 |                     |           |      |               |             |               |
| 331 |                     | Group ID  | Year | UniqueCarrier | ArrTime_AVG | Cancelled_AVG |
| 332 |                     | Group 38  | 1989 | EA            | 1537.789689 | 0.103215      |
| 333 |                     | Group 263 | 2005 | EV            | 1541.068364 | 0.045627      |
| 334 |                     | Group 115 | 2000 | HP            | 1426.261015 | 0.042708      |
| 335 |                     | Group 106 | 2007 | YV            | 1458.986054 | 0.038328      |
| 336 |                     | Group 278 | 2001 | DL            | 1496.292048 | 0.037936      |
| 337 |                     | Group 86  | 2003 | DH            | 1470.618953 | 0.037748      |
| 338 |                     | Group 158 | 2008 | YV            | 1459.219969 | 0.036163      |
| 339 |                     | Group 170 | 2001 | HP            | 1456.840944 | 0.03595       |
| 340 |                     | Group 54  | 1998 | NW            | 1495.768584 | 0.035797      |
| 341 |                     | Group 226 | 2000 | DL            | 1488.152293 | 0.034671      |
| 342 |                     | Group 161 | 2001 | NW            | 1488.901342 | 0.03274       |
| 343 |                     | Group 245 | 1996 | NW            | 1498.360897 | 0.032514      |
| 344 |                     | Group 51  | 2006 | YV            | 1462.74804  | 0.031313      |
| 345 |                     | Group 75  | 2007 | EV            | 1474.859513 | 0.031163      |
| 346 |                     | Group 142 | 2004 | DH            | 1467.669534 | 0.030837      |
| 347 |                     | Group 109 | 2000 | NW            | 1489.586067 | 0.027746      |
| 348 |                     | Group 208 | 2004 | EV            | 1517.992656 | 0.027046      |
| 349 |                     | Group 202 | 2005 | DL            | 1509.912232 | 0.026934      |
| 350 |                     | Group 24  | 2006 | EV            | 1509.957662 | 0.024833      |

Click the down arrow next to Year and sort from largest to smallest. Now we see that Mesa Airlines held the largest on average cancellation percentage in 2008 (.036%). The second least reliable airline in 2008 goes to SkyWest Airlines (OO) with an average cancellation percentage of 2.2% and the third least reliable airline in 2008 is "awarded" to ExpressJet Airlines (EV) with an on average cancellation percentage of 1.8%.

|     | B                   | C         | D    | E             | F           | G             |
|-----|---------------------|-----------|------|---------------|-------------|---------------|
| 329 | <b>Summary Data</b> |           |      |               |             |               |
| 330 |                     |           |      |               |             |               |
| 331 |                     | Group ID  | Year | UniqueCarrier | ArrTime_AVG | Cancelled_AVG |
| 332 |                     | Group 158 | 2008 | YV            | 1459.219969 | 0.036163      |
| 333 |                     | Group 282 | 2008 | OO            | 1460.707471 | 0.021554      |
| 334 |                     | Group 132 | 2008 | EV            | 1458.104271 | 0.017892      |
| 335 |                     | Group 246 | 2008 | B6            | 1443.412562 | 0.015285      |
| 336 |                     | Group 67  | 2008 | DL            | 1495.907908 | 0.015075      |
| 337 |                     | Group 242 | 2008 | AS            | 1480.407939 | 0.014156      |
| 338 |                     | Group 274 | 2008 | OH            | 1491.580041 | 0.009673      |
| 339 |                     | Group 162 | 2008 | FL            | 1498.633192 | 0.008545      |
| 340 |                     | Group 250 | 2008 | NW            | 1500.763541 | 0.008359      |
| 341 |                     | Group 194 | 2008 | MQ            | 1444.519942 | 0.006361      |
| 342 |                     | Group 214 | 2008 | AA            | 1509.266406 | 0.00295       |
| 343 |                     | Group 236 | 2008 | AQ            | 1398.767157 | 0.002571      |
| 344 |                     | Group 55  | 2008 | WN            | 1488.922068 | 0.002303      |
| 345 |                     | Group 137 | 2008 | F9            | 1517.766832 | 0.002132      |
| 346 |                     | Group 244 | 2008 | UA            | 1468.660096 | 0.000844      |
| 347 |                     | Group 269 | 2008 | US            | 1492.253812 | 0.000201      |
| 348 |                     | Group 88  | 2008 | XE            | 1472.404857 | 0.00017       |
| 349 |                     | Group 30  | 2008 | CO            | 1496.828326 | 0.000146      |
| 350 |                     | Group 240 | 2008 | HA            | 1424.607156 | 0.000114      |

Click the down array next to *Cancelled\_AVG* again, but this time sort from *Smallest to Largest*. Then sort by *Year* from *Largest to Smallest*. The table is updated to display the airlines with the smallest on average flight cancellation percentage in 2008: 1<sup>st</sup> place – Hawaiian Airlines, 2<sup>nd</sup> place – Continental Airlines and 3<sup>rd</sup> place – ExpressJet Airlines (1).

|     | B                   | C               | D           | E                    | F                  | G                    |
|-----|---------------------|-----------------|-------------|----------------------|--------------------|----------------------|
| 129 | <b>Summary Data</b> |                 |             |                      |                    |                      |
| 130 |                     |                 |             |                      |                    |                      |
| 131 |                     | <b>Group ID</b> | <b>Year</b> | <b>UniqueCarrier</b> | <b>ArrTime_AVG</b> | <b>Cancelled_AVG</b> |
| 132 |                     | Group 240       | 2008        | HA                   | 1424.607156        | 0.000114             |
| 133 |                     | Group 30        | 2008        | CO                   | 1496.828326        | 0.000146             |
| 134 |                     | Group 88        | 2008        | XE                   | 1472.404857        | 0.00017              |
| 135 |                     | Group 269       | 2008        | US                   | 1492.253812        | 0.000201             |
| 136 |                     | Group 244       | 2008        | UA                   | 1468.660096        | 0.000844             |
| 137 |                     | Group 137       | 2008        | F9                   | 1517.766832        | 0.002132             |
| 138 |                     | Group 55        | 2008        | WN                   | 1488.922068        | 0.002303             |
| 139 |                     | Group 236       | 2008        | AQ                   | 1398.767157        | 0.002571             |
| 140 |                     | Group 214       | 2008        | AA                   | 1509.266406        | 0.00295              |
| 141 |                     | Group 194       | 2008        | MQ                   | 1444.519942        | 0.006361             |
| 142 |                     | Group 250       | 2008        | NW                   | 1500.763541        | 0.008359             |
| 143 |                     | Group 162       | 2008        | FL                   | 1498.633192        | 0.008545             |
| 144 |                     | Group 274       | 2008        | OH                   | 1491.580041        | 0.009673             |
| 145 |                     | Group 242       | 2008        | AS                   | 1480.407939        | 0.014156             |

Using the same steps illustrated here, we could find answers to many other questions, for example:

- What are the yearly flight volume per carrier?
- Which times of day and days of week are most susceptible to departure/arrival delays?
- How many miles per year does each plane by carrier fly?

### **Concluding Remarks**

The ability to sample and summarize large datasets is one that will become more and more important as technology progresses and more and more data is captured. Analytic Solver Data Science's Big Data feature allows users to import these large datasets into Excel allowing business analysts and data scientists the power to build predictive and prescriptive analytic models in their spreadsheets, without the need for complex programming skills. Using Analytic Solver Data Science's Big Data feature, we could easily answer the questions posed in the introduction and more.

# Fitting a model using Feature Selection

---

## What is Feature Selection?

Analytic Solver Data Science's Feature Selection tool gives users the ability to rank and select the most relevant variables for inclusion in a classification or prediction model. In many cases the most accurate models, or the models with the lowest misclassification or residual errors, have benefited from better feature selection, using a combination of human insights and automated methods. Analytic Solver Data Science provides a facility to compute all of the following metrics, described in the literature, to give users information on what features should be included, or excluded, from their models.

- **Correlation-based**
  - Pearson product-moment correlation
  - Spearman rank correlation
  - Kendall concordance
- **Statistical/probabilistic independence metrics**
  - Welch's statistic
  - F statistic
  - Chi-square statistic
- **Information-theoretic metrics**
  - Mutual Information (Information Gain)
  - Gain Ratio
- **Other**
  - Cramer's V
  - Fisher score
  - Gini index

Only some of these metrics can be used in any given application, depending on the characteristics of the input variables (features) and the type of problem. In a supervised setting, if we classify data science problems as follows:

- $\mathbb{R}^n \rightarrow \mathbb{R}$ : real-valued features, prediction (regression) problem
- $\mathbb{R}^n \rightarrow \{0, 1\}$ : real-valued features, binary classification problem
- $\mathbb{R}^n \rightarrow \{1..C\}$ : real-valued features, multi-class classification problem
- $\{1..C\}^n \rightarrow \mathbb{R}^n$ : nominal categorical features, prediction (regression) problem
- $\{1..C\}^n \rightarrow \{0, 1\}$ : nominal categorical features, binary classification problem
- $\{1..C\}^n \rightarrow \{1..C\}$ : nominal categorical features, multi-class classification problem

then we can describe the applicability of the Feature Selection metrics by the following table:

|             | R-R | R-{0,1} | R-{1..C} | {1..C}-R | {1..C}-{0,1} | {1..C}-{1..C} |
|-------------|-----|---------|----------|----------|--------------|---------------|
| Pearson     | N   |         |          |          |              |               |
| Spearman    | N   |         |          |          |              |               |
| Kendall     | N   |         |          |          |              |               |
| Welch's     | D   | N       |          |          |              |               |
| F-Test      | D   | N       | N        |          |              |               |
| Chi-squared | D   | D       | D        | D        | N            | N             |
| Mutual Info | D   | D       | D        | D        | N            | N             |
| Gain Ratio  | D   | D       | D        | D        | N            | N             |
| Fisher      | D   | N       | N        |          |              |               |
| Gini        | D   | N       | N        |          |              |               |

"N" means that metrics can be applied naturally, and "D" means that features and/or the outcome variable must be discretized before applying the particular filter.

As a result, depending on the variables (features) selected and the type of problem chosen in the first dialog, various metrics will be available or disabled in the second dialog.

---

## Feature Selection Example

The goal of this example is three-fold: 1. To use Feature Selection as a tool for exploring relationships between features and the outcome variable, 2. Reducing the dimensionality based on the Feature Selection results and 3. Evaluating the performance of a supervised learning algorithm (a classification algorithm) for different feature subsets.

This example uses the [Heart Failure Clinical Records Dataset](#)<sup>4</sup>, which contains thirteen variables describing 299 patients experiencing heart failure. The [journal article](#) referenced here discusses how the authors analyzed the dataset to first rank the features (variables) by significance and then used the Random Trees machine learning algorithm to fit a model to the dataset. This example attempts to emulate their results.

A description of each variable contained in the dataset appears in the table below.

| VARIABLE | DESCRIPTION    |
|----------|----------------|
| AGE      | Age of patient |

---

<sup>4</sup> Davide Chicco, Giuseppe Jurman: Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone. BMC Medical Informatics and Decision Making 20, 16 (2020). ([link](#))



|                        |                                                                        |
|------------------------|------------------------------------------------------------------------|
| ANAEMIA                | Decrease of red blood cells or hemoglobin (boolean)                    |
| CREATINE_PHOSPHOKINASE | Level of the CPK enzyme in the blood (mcg/L)                           |
| DIABETES               | If the patient has diabetes (boolean)                                  |
| EJECTION_FRACTION      | Percentage of blood leaving the heart at each contraction (percentage) |
| HIGH_BLOOD_PRESSURE    | If the patient has hypertension (boolean)                              |
| PLATELETS              | Platelets in the blood (kiloplatelets/mL)                              |
| SERUM_CREATININE       | Level of serum creatinine in the blood (mg/dL)                         |
| SERUM_SODIUM           | Level of serum sodium in the blood (mEq/L)                             |
| SEX                    | Woman (0) or man (1)                                                   |
| SMOKING                | If the patient smokes or not (boolean)                                 |
| TIME                   | Follow-up period (days)                                                |
| DEATH_EVENT            | If the patient deceased during the follow-up period (boolean)          |

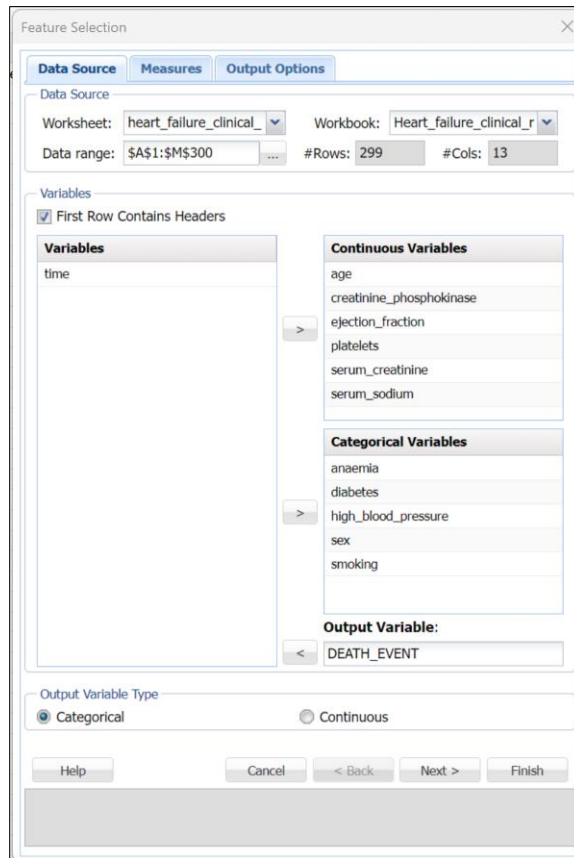
To open the example dataset, click **Help – Example Models – Forecasting/Data Science Examples – Heart Failure Clinical Records**.

Select a cell within the data (say A2), then click **Explore – Feature Selection** to bring up the first dialog.

- Select all the following variables as Continuous Variables: age, creatinine\_phosphokinase, ejection\_fraction, platelets, serum\_creatinine and serum\_sodium.
- Select the following variables as Categorical Variables: anaemia, diabetes, high\_blood\_pressure, sex and smoking.
- Select DEATH\_EVENT as the Output Variable
- Confirm that Categorical is selected for Output Variable Type. This setting denotes that the Output Variable is a categorical variable. If the number of unique values in the *Output variable* is greater than 10, then *Continuous* will be selected by default. However, at any time the User may override the default choice based on his or her own knowledge of the variable.

This analysis omits the time variable. The Feature Selection dialog should look similar to the screenshot below.

Figure 1: Feature Selection Data Source dialog



Click the **Measures** tab or click **Next** to open the Measures dialog.

Since we have continuous variables, *Discretize predictors* is enabled. When this option is selected, Analytic Solver Data Science will transform continuous variables into discrete, categorical data in order to be able to calculate statistics, as shown in the table in the Introduction to this chapter.

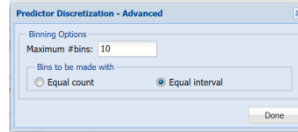
This dataset contains both continuous (or real-valued) features and categorical features which puts this dataset into the following category.

$\mathbb{R}^n \rightarrow \{0, 1\}$ : real-valued features, binary classification problem

As a result, if interested in evaluating the relevance of features according to the Chi-Squared Test or measures available in the Information Theory group (Mutual Information and Gain ratio), the variables must first be discretized.

Select **Discretize predictors**, then click **Advanced**. Leave the defaults of **10** for *Maximum # bins* and **Equal Interval** for *Bins to be made with*. Analytic Solver Data Science will create 10 bins and will assign records to the bins based on if the variable's value falls in the interval of the bin. This will be performed for each of the Continuous Variables.

Figure 2: Predictor Discretization - Advanced



Note: Discretize output variable is disabled because our output variable, DEATH\_EVENT, is already a categorical nominal variable. If we had no Continuous Variables and all Categorical Variables, *Discretize predictors* would be disabled.

Select **Chi-squared** and **Cramer's V** under *Chi-Squared Test*. The Chi-squared test statistic is used to assess the statistical independence of two events. When applied to Feature Selection, it is used as a test of independence to assess whether the assigned class is independent of a particular variable. The minimum value for this statistic is 0. The higher the Chi-Squared statistic, the more independent the variable.

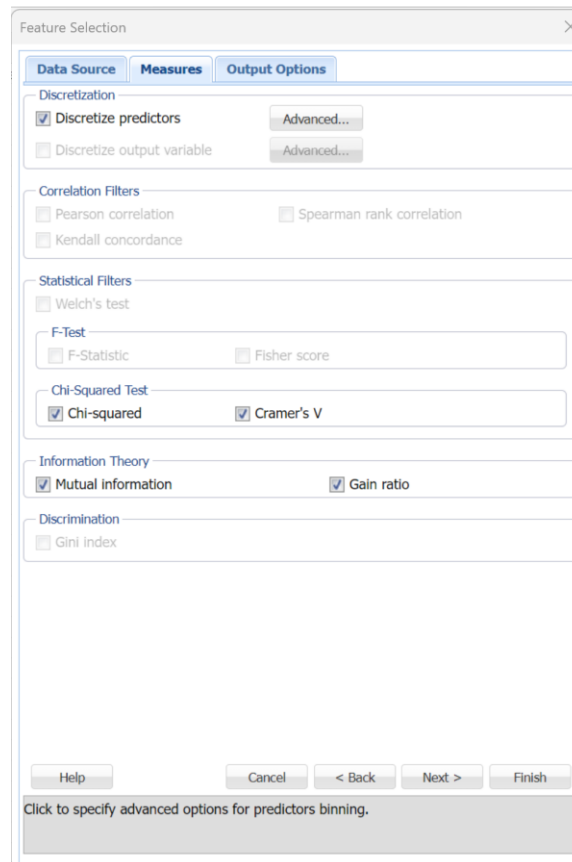
Cramer's V is a variation of the Chi-Squared statistic that also measures the association between two discrete nominal variables. This statistic ranges from 0 to 1 with 0 indicating no association between the two variables and 1 indicating complete association (the two variables are equal).

Select *Mutual information* and *Gain ratio* within the *Information Theory* frame. Mutual information is the degree of a variables' mutual dependence or the amount of uncertainty in variable 1 that can be reduced by incorporating knowledge about variable 2. Mutual Information is non-negative and is equal to zero if the two variables are statistically independent. Also, it is always less than the entropy (amount of information contained) in each individual variable.

The *Gain Ratio*, ranging from 0 and 1, is defined as the mutual information (or information gain) normalized by the feature entropy. This normalization helps address the problem of overemphasizing features with many values but the normalization results in an overestimate of the relevance of features with low entropy. It is a good practice to consider both mutual information and gain ratio for deciding on feature rankings. The larger the gain ratio, the larger the evidence for the feature to be relevant in a classification model.

For more information on the remaining options on this dialog, see the Using Feature Selection found within the Data Science Reference Guide.

Figure 3: Feature Selection Measures dialog



Click the **Output Options** tab or click **Next** to open the Output Options dialog. *Table of all produced measures* is selected by default. When this option is selected, Analytic Solver Data Science will produce a report containing all measures selected on the Measures tab.

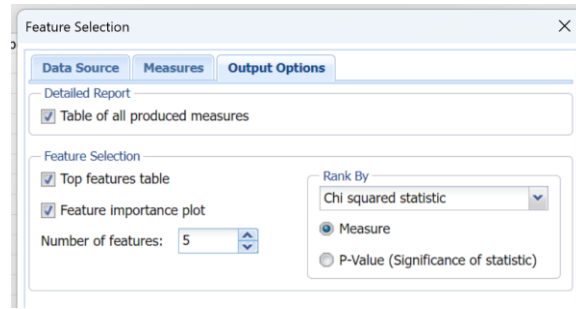
**Top Features table** is selected by default. This option produces a report containing only top variables as indicated by the Number of features edit box.

Select **Feature importance plot**. This option produces a graphical representation of variable importance based on the measure selected in the *Rank By* drop down menu.

Enter **5** for *Number of features*. Analytic Solver Data Science will display the top 5 most important or most relevant features (variables) as ranked by the statistic displayed in the *Rank By* drop down menu.

Keep *Chi squared statistic* selected for the *Rank By* option. Analytic Solver Data Science will display all measures and rank them by the statistic chosen in this drop down menu.

Figure 4: Feature Selection Output Options dialog



Click **Finish**.

Two worksheets are inserted to the right of the heart\_failure\_clinical\_records worksheet: FS\_Output and FS\_Top\_Features.

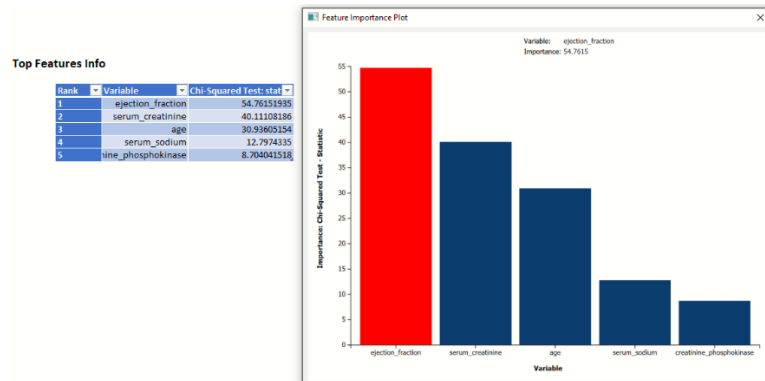
Click the FS\_Top\_Features tab.

In the Data Science Cloud app, click the Charts icon on the Ribbon to open the Charts dialog, then select **FS\_Top\_Features** for *Worksheet* and **Feature Importance Chart** for *Chart*.

The Feature Importance Plot ranks the variables by most important or relevant according to the selected measure. In this example, we see that the ejection\_fraction, serum\_creatinine, age, serum\_sodium and creatinine\_phosphokinase are the top five most important or relevant variables according to the Chi-Squared statistic. It's beneficial to examine the Feature Selection Importance Plot in order to quickly identify the largest drops or "elbows" in feature relevancy (importance) and select the optimal number of variables for a given classification or regression model.

Note: We could have limited the number of variables displayed on the plot to a specified number of variables (or features) by selecting Number of features and then specifying the number of desired variables. This is useful when the number of input variables is large or we are particularly interested in a specific number of highly – ranked features.

Figure 5: Top Features Plot and Table



Run your mouse over each bar in the graph to see the Variable name and Importance factor, in this case Chi-Square, in the top of the dialog.

Click the **X** in the upper right hand corner to close the dialog, then click *FS\_Output* tab to open the Feature Selection report.

Figure 6: Feature Selection: Statistics Table

| Variable                 | Chi2: stat | Chi2: p-value | Chi2: Cramer's V | Mutual Information | Gain Ratio  |
|--------------------------|------------|---------------|------------------|--------------------|-------------|
| age                      | 30.9360515 | 0.000303603   | 0.321659844      | 0.074279228        | 0.024162533 |
| creatinine_phosphokinase | 8.70404152 | 0.27460876    | 0.170618014      | 0.023222979        | 0.024270684 |
| ejection_fraction        | 54.7615193 | 1.35227E-08   | 0.427958987      | 0.130937229        | 0.051493108 |
| platelets                | 7.77686143 | 0.556780346   | 0.161274828      | 0.021613206        | 0.009986181 |
| serum_creatinine         | 40.1110819 | 3.05463E-06   | 0.36626599       | 0.096222047        | 0.074732546 |
| serum_sodium             | 12.7974335 | 0.11901195    | 0.206883496      | 0.030659794        | 0.013665886 |
| anaemia                  | 1.31312606 | 0.251829444   | 0.066270098      | 0.003156983        | 0.00320053  |
| diabetes                 | 0.00112866 | 0.973199636   | 0.001942883      | 2.72339E-06        | 2.77744E-06 |
| high_blood_pressure      | 1.88268052 | 0.170029802   | 0.079351058      | 0.004494526        | 0.004806418 |
| sex                      | 0.0055707  | 0.940503436   | 0.004316376      | 1.34304E-05        | 1.43623E-05 |
| smoking                  | 0.04764385 | 0.827215074   | 0.012623153      | 0.000115234        | 0.000127255 |

The Detailed Feature Selection Report displays each computed metric selected on the Measures tab: Chi-squared statistic, Chi-squared P-Value, Cramer’s V, Mutual Information, and Gain Ratio.

### Chi2: Statistic and p-value

Click the down arrow next to Chi2: p-value to sort the table according to this statistic going from smallest p-value to largest.

Figure 7: Statistics sorted by Chi2:p-value

| Variable                 | Chi2: stat | Chi2: p-value |
|--------------------------|------------|---------------|
| ejection_fraction        | 54.7615193 | 1.35227E-08   |
| serum_creatinine         | 40.1110819 | 3.05463E-06   |
| age                      | 30.9360515 | 0.000303603   |
| serum_sodium             | 12.7974335 | 0.11901195    |
| high_blood_pressure      | 1.88268052 | 0.170029802   |
| anaemia                  | 1.31312606 | 0.251829444   |
| creatinine_phosphokinase | 8.70404152 | 0.27460876    |
| platelets                | 7.77686143 | 0.556780346   |
| smoking                  | 0.04764385 | 0.827215074   |
| sex                      | 0.0055707  | 0.940503436   |
| diabetes                 | 0.00112866 | 0.973199636   |

According to the Chi-squared test, ejection\_fraction, serum\_creatinine and age are the 3 most relevant variables for predicting the outcome of a patient in heart failure.

### Chi2: Cramer's V

Recall that the Cramer's V statistic ranges from 0 to 1 with 0 indicating no association between the two variables and 1 indicating complete association (the two variables are equal). Sort the Cramer's V statistic from largest to smallest.

Figure 8: Chi2: Cramer's V Statistic

| Variable                 | Chi2: Cramer's V |
|--------------------------|------------------|
| ejection_fraction        | 0.427958987      |
| serum_creatinine         | 0.36626599       |
| age                      | 0.321659844      |
| serum_sodium             | 0.206883496      |
| creatinine_phosphokinase | 0.170618014      |
| platelets                | 0.161274828      |
| high_blood_pressure      | 0.079351058      |
| anaemia                  | 0.066270098      |
| smoking                  | 0.012623153      |
| sex                      | 0.004316376      |
| diabetes                 | 0.001942883      |

Again, this statistic ranks the same four variables, ejection\_fraction, serum\_creatinine, age and serum\_sodium, as the Chi<sup>2</sup> statistic.

## Mutual Information

Sort the Mutual Information column by largest to smallest value. This statistic measures how much information the presence/absence of a term contributes to making the correct classification decision.<sup>5</sup> The closer the value to 1, the more contribution the feature provides.

Figure 9: Mutual Information Statistic

### Feature Selection: Statistics

| Variable                 | Mutual Information | Ge |
|--------------------------|--------------------|----|
| ejection_fraction        | 0.130937229        | 0  |
| serum_creatinine         | 0.096222047        | 0  |
| age                      | 0.074279228        | 0  |
| serum_sodium             | 0.030659794        | 0  |
| creatinine_phosphokinase | 0.023222979        | 0  |
| platelets                | 0.021613206        | 0  |
| high_blood_pressure      | 0.004494526        | 0  |
| anaemia                  | 0.003156983        | 0  |
| smoking                  | 0.000115234        | 0  |
| sex                      | 1.34304E-05        | :  |
| diabetes                 | 2.72339E-06        | :  |

When compared to the Chi2 and Cramer's V statistic, the top four most significant variables calculated for Mutual Information are the same: ejection\_fraction, serum\_creatinine, age, and serum\_sodium.

## Gain Ratio

Finally, sort the Gain Ratio from largest to smallest. (Recall that the larger the gain ratio value, the larger the evidence for the feature to be relevant in the classification model.)

Figure 10: Gain Ratio

### Feature Selection: Statistics

| Variable                 | Gain Ratio | Ge |
|--------------------------|------------|----|
| serum_creatinine         | 0.07473255 |    |
| ejection_fraction        | 0.05149311 |    |
| creatinine_phosphokinase | 0.02427068 |    |
| age                      | 0.02416253 |    |
| serum_sodium             | 0.01366589 |    |
| platelets                | 0.00998618 |    |
| high_blood_pressure      | 0.00480642 |    |
| anaemia                  | 0.00320053 |    |
| smoking                  | 0.00012725 |    |
| sex                      | 1.4362E-05 |    |
| diabetes                 | 2.7774E-06 |    |

While this statistic's rankings differ from the first 4 statistic's rankings, ejection\_fraction, age and serum\_creatinine are still ranked in the top four positions.

The Feature Selection tool has allowed us to quickly explore and learn about our data. We now have a pretty good idea of which variables are the most relevant or most important to our classification or prediction model, how our variables relate to each other and to the output variable, and which data attributes would be worth extra time and money in future data collection. Interestingly, for this example, most of our ranking statistics have agreed (mostly) on the most important or relevant features with strong evidence. We computed and

<sup>5</sup> <https://nlp.stanford.edu/IR-book/html/htmledition/mutual-information-1.html>

examined various metrics and statistics and for some (where p-values can be computed) we've seen a statistical evidence that the test of interest succeeded with definitive conclusion. In this example, we've observed that several variable (or features) were consistently ranked in the top 3-4 most important variables by most of the measures produced by Analytic Solver Data Science's Feature Selection tool. However, this will not always be the case. On some datasets you will find that the ranking statistics and metrics compete on rankings. In cases such as these, further analysis may be required.

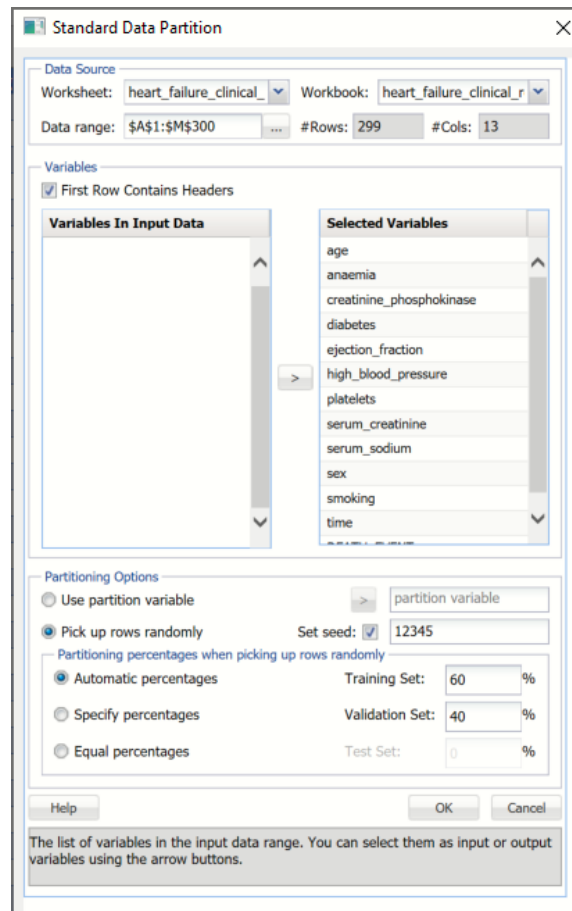
### **Fitting the Model**

To further emulate the results of the journal article discussed in the Feature Selection Example, the Random Trees Ensemble Classification Methods will be used to investigate if a machine learning algorithm can predict a patient's survival using the top two or three ranked features as found by the Feature Selection tool.

First, click **Partition – Standard Partition** to partition the dataset into Training, Validation and Test Sets using the default percentages of 60% allocated to the Training Set and 40% allocate to the Validation Set.

For more information on the remaining options on this dialog, see the Standard Partitioning chapter within the Data Science Reference Guide.

Figure 11: Standard Data Partition dialog



Then click **OK** to create the two partitions. A new worksheet *STDPartition* is inserted to the right of the dataset. The number of records allocated to the



Training partition is 179 and the number of records allocated to the Validation partition is 120.

Figure 12: Standard Data Partitioning results

| Record ID | age | anaemia | creatinine_phosphokinase | diabetes | ejection_fraction | high_blood_pressure | platelets | serum_creatinine | serum_sodium | sex | smoking | time |
|-----------|-----|---------|--------------------------|----------|-------------------|---------------------|-----------|------------------|--------------|-----|---------|------|
| Record 1  | 75  | 0       | 582                      | 0        | 20                | 1                   | 265000    | 1.9              | 130          | 1   | 0       |      |
| Record 5  | 90  | 1       | 47                       | 0        | 40                | 1                   | 204000    | 2.1              | 132          | 1   | 1       |      |
| Record 8  | 80  | 1       | 123                      | 0        | 35                | 1                   | 386000    | 9.4              | 133          | 1   | 1       |      |
| Record 15 | 45  | 0       | 582                      | 0        | 14                | 0                   | 166000    | 0.8              | 127          | 1   | 0       |      |
| Record 17 | 65  | 0       | 146                      | 1        | 38                | 0                   | 146000    | 1.0              | 144          | 1   | 1       |      |

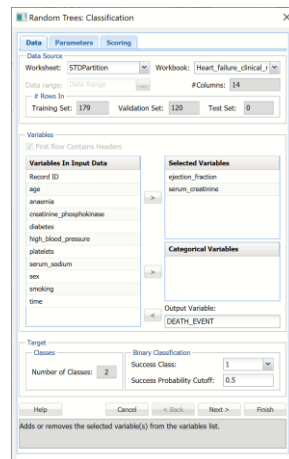
The first time that the model is fit, only two features (ejection\_fraction and serum\_creatinine) will be utilized.

Click **Classify – Ensemble – Random Trees** to open the *Random Trees: Classification* dialog.

- Select the two Variables from *Variables In Input Data* (ejection\_fraction and serum\_creatinine) and click the right pointing arrow to the left of *Selected Variables* to add these two variables to the model. Then take similar steps to select DEATH\_EVENT as the *Output Variable*.
- Leave Success Class as "1" and Success Probability Cutoff at 0.5 under Binary Classification.

The Random Trees: Classification dialog should be similar to the one pictured in the Figure 10 below.

Figure 13: Random Trees: Classification dialog with Selection Variables (serum\_creatinine and ejection\_fraction) and Output Variables (DEATH\_EVENT) selected.



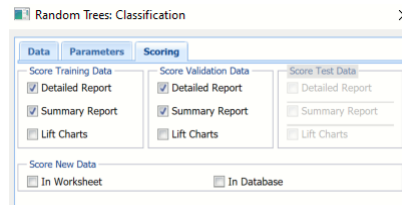
Click the **Scoring** tab to advance to the *Random Trees: Classification Scoring* tab.

For more information on Random Trees parameters, see the Random Trees Classification Options section below.

Summary Report is selected by default. Select **Detailed Report** for both *Score Training Data* and *Score Validation Data* and then click **Finish**.

For more information on the remaining options on any of the Random Trees: Classification dialogs, see the Random Trees: Classification Options section within the Data Science Reference Guide.

Figure 14: Random Trees: Classification dialog with output choices selected



Four worksheets are inserted to the right of the STDPartition tab: *CRandTrees\_Output*, *CRandTrees\_TrainingScore*, *CRandTrees\_ValidationScore* and *CRandTrees\_Stored*.

- *CRandTrees\_Output* reports the input data, output data, and parameter settings.
- *CRandTrees\_TrainingScore* reports the confusion matrix, calculated metrics and the actual classification by row for the training partition.
- *CRandTrees\_ValidationScore* reports the confusion matrix, calculated metrics and the actual classification by row for the validation partition.
- *CRandTrees\_Stored* contains the stored model which can be used to apply the fitted model to new data. See the Scoring chapter within the Analytic Solver Data Science User Guide for an example of scoring new data using the stored model.

Click **CRandTrees\_TrainingScore** to view the Classification Summary for the Training partition.

Figure 15: Training: Classification Summary

| Training: Classification Summary        |             |                         |             |             |
|-----------------------------------------|-------------|-------------------------|-------------|-------------|
| <b>Confusion Matrix</b>                 |             |                         |             |             |
| Actual\Predicted                        | 0           | 1                       |             |             |
| 0                                       | 103         | 14                      |             |             |
| 1                                       | 17          | 45                      |             |             |
| <b>Error Report</b>                     |             |                         |             |             |
| Class                                   | # Cases     | # Errors                | % Error     |             |
| 0                                       | 117         | 14                      | 11.96581197 |             |
| 1                                       | 62          | 17                      | 27.41935484 |             |
| Overall                                 | 179         | 31                      | 17.31843575 |             |
| <b>Metrics</b>                          |             |                         |             |             |
| Metric                                  | Value       |                         |             |             |
| Accuracy (#correct)                     | 148         |                         |             |             |
| Accuracy (%correct)                     | 82.68156425 |                         |             |             |
| Specificity                             | 0.88034188  |                         |             |             |
| Sensitivity (Recall)                    | 0.725806452 |                         |             |             |
| Precision                               | 0.762711864 |                         |             |             |
| F1 score                                | 0.743801653 |                         |             |             |
| Success Class                           | 1           |                         |             |             |
| Success Probability                     | 0.5         |                         |             |             |
| <b>Training: Classification Details</b> |             |                         |             |             |
| Record ID                               | DEATH_EVENT | Prediction: DEATH_EVENT | PostProb: 0 | PostProb: 1 |
| Record 1                                | 1           | 1                       | 0           | 1           |
| Record 5                                | 1           | 1                       | 0           | 1           |
| Record 8                                | 1           | 0                       | 0.8         | 0.2         |
| Record 15                               | 0           | 1                       | 0.5         | 0.5         |
| Record 18                               | 1           | 1                       | 0.1         | 0.9         |

The overall error for the training partition was 17.32% with 14 surviving patients reported as deceased and 17 deceased patients reported as survivors.

- Accuracy: 82.68% -- Accuracy refers to the ability of the classifier to predict a class label correctly.
- Specificity: 0.88 – (True Negative)/(True Negative + False Positives)

Specificity is defined as the proportion of negative classifications that were actually negative, or the fraction of survivors that actually survived. In this model, 103 actual surviving patients were classified correctly as survivors.

There were 14 false positives or 14 actual survivors classified incorrectly as deceased.

- Sensitivity or Recall:  $0.726 = (\text{True Positive})/(\text{True Positive} + \text{False Negative})$

Sensitivity is defined as the proportion of positive cases there were classified correctly as positive, or the proportion of actually deceased patients there were classified as deceased. In this model, 45 actual deceased patients were correctly classified as deceased. There were 17 false negatives or 17 actual deceased patients were incorrectly classified as survivors.

Note: Since the object of this model is to correctly classify which patients will succumb to heart failure, this is an important statistic as it is very important for a physician to be able to accurately predict which patients require mitigation.

- Precision:  $0.763 = (\text{True Positives})/(\text{True Positives} + \text{False Positives})$

Precision is defined as the proportion of positive results that are true positive. In this model, 45 actual deceased patients were classified correctly as deceased. There were 14 false positives or 14 actual survivors classified incorrectly as deceased.

- F-1 Score:  $0.743 = 2 \times (\text{Precision} * \text{Sensitivity})/(\text{Precision} + \text{Sensitivity})$

The F-1 Score provides a statistic to balance between Precision and Sensitivity, especially if an uneven class distribution exists, as in this example, (103 survivors vs 45 deceased). The closer the F-1 score is to 1 (the upper bound) the better the precision and recall.

- Success Class and Success Probability simply reports the settings for these two values as input on the Random Trees: Classification, Data tab.
- View individual records and their classifications beneath Training: Classification Details.

Click **the CRandTrees\_ValidationScore** tab to view the Summary Results for the Validation partition.

Figure 16: Validation: Classification Summary

|    | B                                  | C           | D                       | E           | F           | G |
|----|------------------------------------|-------------|-------------------------|-------------|-------------|---|
| 10 | Validation: Classification Summary |             |                         |             |             |   |
| 11 |                                    |             |                         |             |             |   |
| 12 | Confusion Matrix                   |             |                         |             |             |   |
| 13 | Actual\Predicted                   | 0           | 1                       |             |             |   |
| 14 | 0                                  | 67          | 19                      |             |             |   |
| 15 | 1                                  | 10          | 24                      |             |             |   |
| 16 |                                    |             |                         |             |             |   |
| 17 | Error Report                       |             |                         |             |             |   |
| 18 | Class                              | # Cases     | # Errors                | % Error     |             |   |
| 19 | 0                                  | 86          | 19                      | 22.09302326 |             |   |
| 20 | 1                                  | 34          | 10                      | 29.41176471 |             |   |
| 21 | Overall                            | 120         | 29                      | 24.16666667 |             |   |
| 22 |                                    |             |                         |             |             |   |
| 23 | Metrics                            |             |                         |             |             |   |
| 24 | Metric                             | Value       |                         |             |             |   |
| 25 | Accuracy (%correct)                | 91          |                         |             |             |   |
| 26 | Accuracy (%correct)                | 75.83333333 |                         |             |             |   |
| 27 | Specificity                        | 0.779069767 |                         |             |             |   |
| 28 | Sensitivity (Recall)               | 0.705882353 |                         |             |             |   |
| 29 | Precision                          | 0.558139535 |                         |             |             |   |
| 30 | F1 score                           | 0.623376623 |                         |             |             |   |
| 31 | Success Class                      | 1           |                         |             |             |   |
| 32 | Success Probability                | 0.5         |                         |             |             |   |
| 33 |                                    |             |                         |             |             |   |
| 34 | Validation: Classification Details |             |                         |             |             |   |
| 35 |                                    |             |                         |             |             |   |
| 36 | Record ID                          | DEATH_EVENT | Prediction: DEATH_EVENT | PostProb: 0 | PostProb: 1 |   |
| 37 | Record 104                         | 0           | 1                       | 0.5         | 0.5         |   |
| 38 | Record 163                         | 0           | 0                       | 1           | 0           |   |
| 39 | Record 290                         | 0           | 0                       | 1           | 0           |   |
| 40 | Record 207                         | 0           | 0                       | 1           | 0           |   |
| 41 | Record 126                         | 0           | 0                       | 0.7         | 0.3         |   |

The overall error for the validation partition was 24.17 with 19 false positives (surviving patients reported as deceased) and 10 false negatives (deceased patients reported as survivors).

Note the following metrics:

- Accuracy: 75.83
- Specificity: 0.779
- Sensitivity or Recall: 0.706
- Precision: 0.558
- F1 Score: 0.623

These steps were performed multiple times while adding additional Selected Variables according to the variable's importance or significance found by Feature Selection. The results are summarized in the table below.

The lowest Overall Error in the Validation Partition for any of the variable combinations occurs when just two variables, ejection\_fraction and serum\_creatinine, are present in the fitted model. In addition, this fitted model also exhibits the highest Accuracy, Sensitivity, Precision and F1 Score metrics. These results suggest that by obtaining these two measurements for a patient, a physician can determine whether the patient should undergo some type of mitigation for their heart failure diagnosis.

| Variables                                                | Training Partition |                      |             |                      |           |          | Validation Partition |          |             |                      |           |          |
|----------------------------------------------------------|--------------------|----------------------|-------------|----------------------|-----------|----------|----------------------|----------|-------------|----------------------|-----------|----------|
|                                                          | Overall Error      | Accuracy (% Correct) | Specificity | Sensitivity (Recall) | Precision | F1 Score | Overall Error        | Accuracy | Specificity | Sensitivity (Recall) | Precision | F1 Score |
| ejection_fraction, serum_creatinine                      | 17.318             | 82.682               | 0.880       | 0.726                | 0.763     | 0.744    | 24.167               | 75.833   | 0.779       | 0.706                | 0.558     | 0.623    |
| + age                                                    | 3.352              | 96.648               | 0.966       | 0.968                | 0.938     | 0.952    | 30.000               | 70.000   | 0.744       | 0.588                | 0.476     | 0.526    |
| + serum_sodium                                           | 3.911              | 96.089               | 0.966       | 0.952                | 0.937     | 0.944    | 28.333               | 71.667   | 0.767       | 0.588                | 0.500     | 0.541    |
| + high_blood_pressure<br>(added as categorical variable) | 2.793              | 97.207               | 0.966       | 0.984                | 0.938     | 0.961    | 29.167               | 70.833   | 0.756       | 0.588                | 0.488     | 0.533    |
| + anaemia<br>(added as categorical variable)             | 5.587              | 94.413               | 0.966       | 0.903                | 0.933     | 0.918    | 27.500               | 72.500   | 0.814       | 0.500                | 0.515     | 0.507    |
| + creatinine_phosphokinase                               | 2.235              | 97.765               | 0.974       | 0.984                | 0.953     | 0.968    | 30.833               | 69.167   | 0.744       | 0.559                | 0.463     | 0.507    |
| + platelets                                              | 2.235              | 97.765               | 0.974       | 0.984                | 0.953     | 0.968    | 30.833               | 69.167   | 0.744       | 0.559                | 0.463     | 0.507    |
| + smoking<br>(added as categorical variable)             | 3.352              | 96.648               | 0.983       | 0.935                | 0.967     | 0.951    | 29.167               | 70.833   | 0.756       | 0.588                | 0.488     | 0.533    |
| + sex<br>(added as categorical variable)                 | 3.352              | 96.648               | 0.974       | 0.952                | 0.952     | 0.952    | 29.167               | 70.833   | 0.791       | 0.500                | 0.486     | 0.493    |
| + diabetes<br>(added as categorical variable)            | 3.911              | 96.089               | 0.974       | 0.935                | 0.951     | 0.943    | 30.833               | 69.167   | 0.767       | 0.500                | 0.459     | 0.479    |

# Text Mining

---

## Introduction

Text mining is the practice of automated analysis of one document or a collection of documents (corpus) and extracting non-trivial information from it. Also, Text Mining usually involves the process of transforming unstructured textual data into structured representation by analyzing the patterns derived from text. The results can be analyzed to discover interesting knowledge, some of which would only be found by a human carefully reading and analyzing the text. Typical widely-used tasks of Text Mining include but are not limited to Automatic Text Classification/Categorization, Topic Extraction, Concept Extraction, Documents/Terms Clustering, Sentiment Analysis, Frequency-based Analysis and many more. Some of these tasks could not be completed by a human, which makes Text Mining a particularly useful and applicable tool in modern Data Science. Analytic Solver Data Science's Text Miner takes an integrated approach to text mining as it does not totally separate analysis of unstructured data from traditional data science techniques applicable for structured information. While Analytic Solver Data Science is a very powerful tool for analyzing text only, it also offers automated treatment of mixed data, i.e. combination of multiple unstructured and structured fields. This is a particularly useful feature that has many real-world applications, such as analyzing maintenance reports, evaluation forms, insurance claims, etc. Text Miner uses the "bag of words" model – the simplified representation of text, where the precise grammatical structure of text and exact word order is disregarded. Instead, syntactic, frequency-based information is preserved and is used for text representation. Although such assumptions might be harmful for some specific applications of Natural Language Processing (NLP), it has been proven to work very well for applications such as Text Categorization, Concept Extraction and others, which are the particular areas addressed by Text Miner's capabilities. It has been shown in many theoretical/empirical studies that syntactic similarity often implies semantic similarity. One way to access syntactic relationships is to represent text in terms of Generalized Vector Space Model (GVSP). Advantage of such representation is a meaningful mapping of text to the numeric space, the disadvantage is that some semantic elements, e.g. order of words, are lost (recall the bag-of-words assumption).

Input to Text Miner (the Text Mining tool within Analytic Solver Data Science) could be of two main types – few relatively large documents (e.g. several books) or relatively large number of smaller documents (e.g. collection of emails, news articles, product reviews, comments, tweets, Facebook posts, etc.). While Text Miner is capable of analyzing large text documents, it is particularly effective for large corpuses of relatively small documents. Obviously, this functionality has limitless number of applications – for instance, email spam detection, topic extraction in articles, automatic rerouting of correspondence, sentiment analysis of product reviews and many more.

The input for text mining is a dataset on a worksheet, with at least one column that contains free-form text (or file paths to documents in a file system containing free-form text), and, optionally, other columns that contain traditional structured data. In the first tab of the Text Mining dialog, the user selects the text variable(s), and the other variable(s) to be processed.

The output for the text mining is a set of reports that contain general explorative information about the collection of documents and structured representations of text (free-form text columns are expanded to a set of new columns with numeric representation. The new columns will each correspond to either (i) a single term (word) found in the “corpus” of documents, or, if requested, (ii) a concept extracted from the corpus through Latent Semantic Indexing (LSI, also called LSA or Latent Semantic Analysis). Each concept represents an automatically derived complex combination of terms/words that have been identified to be related to a particular topic in the corpus of documents. The structural representation of text can serve as an input to any traditional data science techniques available in Text Miner – unsupervised/supervised, affinity, visualization techniques, etc. In addition, Text Miner also presents a visual representation of Text Mining results to allow the user to interactively explore the information, which otherwise would be extremely hard to analyze manually. Typical visualizations that aid in understanding of Text Mining outputs and that are produced by Text Miner are:

- Zipf plot – for visual/interactive exploration of frequency-based information extracted by Text Miner
- Scree Plot, Term-Concept and Document-Concept 2D scatter plots – for visual/interactive exploration of Concept Extraction results

If you are interested in visualizing specific parts of Text Mining analysis outputs, Text Miner provides rich capabilities for charting – the functionality that can be used to explore Text Mining results and supplement standard charts discussed above.

In the example below, you will learn how to use Text Miner in Analytic Solver Data Science to process/analyze approximately 1000 text files and use the results for automatic topic categorization. This will be achieved by using structured representation of text presented to Logistic Regression for building the model for classification.

---

## Text Mining Example

This example uses the text files within the Text Mining Example Documents.zip archive file to illustrate how to use Analytic Solver Data Science’s Text Mining tool. These documents were selected from the well-known text dataset (downloadable from <http://www.cs.cmu.edu/afs/cs/project/theo-20/www/data/news20.html>) which consists of 20,000 messages, collected from 20 different internet newsgroups. We selected about 1,200 of these messages that were posted to two interest groups, Autos and Electronics (about 500 documents from each).

Note: Analytic Solver Cloud does not currently support importing from a file folder.

### Importing from a File Folder

The Text Mining Example Documents.zip archive file is located at C:\ProgramData\Frontline Systems\Datasets. Unzip the contents of this file to a location of your choice. Four folders will be created beneath Text Mining Example Documents: Autos, Electronics, Additional Autos and Additional Electronics. One thousand, two hundred short text files will be extracted to the location chosen. This example is based on the text dataset at <http://www.cs.cmu.edu/afs/cs/project/theo-20/www/data/news20.html>, which consists of 20,000 messages, collected from 20 different netnews newsgroups.

We selected about 1,200 of these messages that were posted to two interest groups, for Autos and Electronics (about 50% in each).

Select **Get Data – File Folder** to open the *Import From File System* dialog.

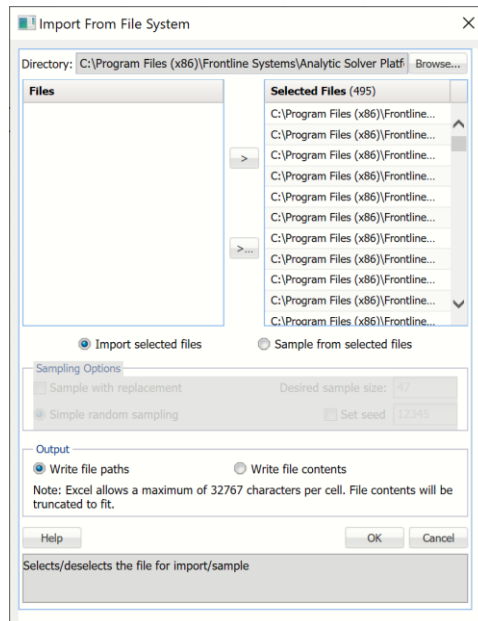
At the top of the dialog, click **Browse...** to navigate to the *Autos* subfolder (C:\ProgramData\Frontline Systems\Datasets\Text Mining Example Documents\Autos). Set the *File Type* to All Files (\*.\*), then select all files in the folder and click the **Open** button. The files will appear in the left list box under *Files*. Click the >> button to move the files from the *Files* list box to the *Selected Files* list box. Now repeat these steps for the *Electronics* subfolder. When these steps are completed, 985 files will appear under *Selected Files*.

Select **Sample from selected files** to enable the *Sampling Options*. Text Miner will perform sampling from the files in the *Selected Files* field. Enter **300** for *Desired sample size* while leaving the default settings for *Simple random sampling* and *Set Seed*.

*Note: If you are using the educational version of Analytic Solver Data Science, enter "100" for Desired Sample Size. This is the upper limit for the number of files supported when sampling from a file system when using Analytic Solver Data Science. For a complete list of the capabilities of Analytic Solver Data Science and Analytic Solver Data Science for Education, click [here](#).*

Text Miner will select 300 files using Simple random sampling with a seed value of 12345. Under *Output*, leave the default setting of *Write file paths*. Rather than writing out the file contents into the report, Text Miner will include the file paths.

*Note: Currently, Analytic Solver Data Science only supports the import of delimited text files. A delimited text file is one in which data values are separated by a character such as quotation marks, commas or tabs. These characters define a beginning and end of a string of text.*



Click **OK**. The output *XLM\_SampleFiles* will be inserted into the Data Science task pane, with contents similar to that shown on the next page.

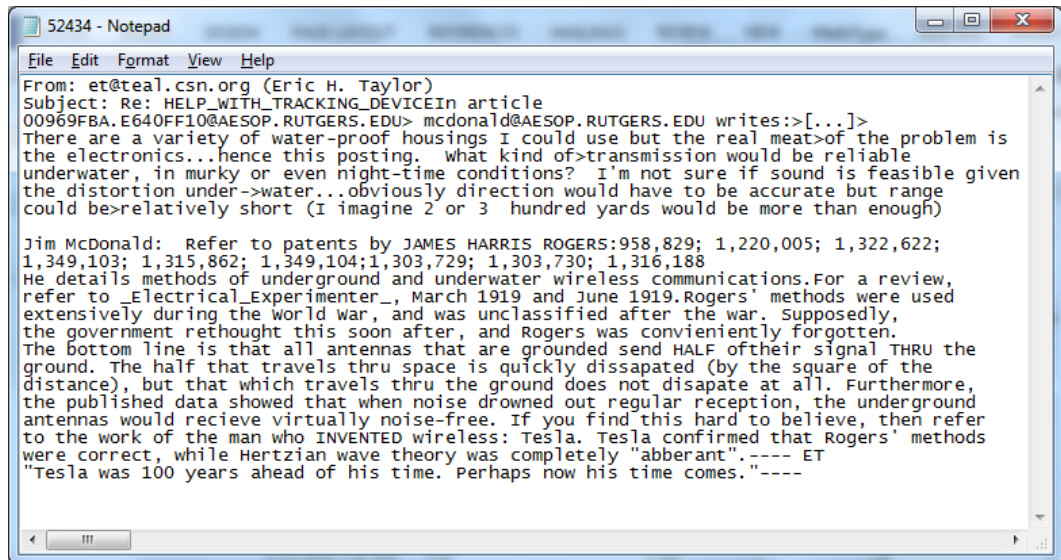
The *Data* portion of the report displays the selections we made on the *Import From File System* dialog. Here we see the path of the directories, the number of

files written, our choice to write the paths or contents (*File Paths*), the sampling method, the desired sample size, the actual size of the sample, and the seed value (12345).

Underneath the *Data* portion are paths to the 300 text files in random order that were sampled by Analytic Solver Data Science. If *Write file contents* had been selected, rather than *Write file paths*, the report would contain the RowID, File Path, and the first 32,767 characters present in the document.

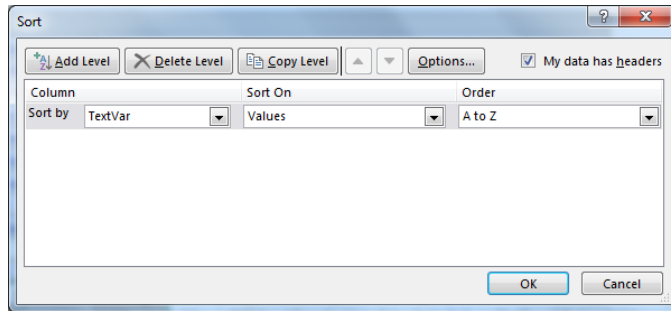
| A  | B                          | C                                          | D | E | F | G | H | I |
|----|----------------------------|--------------------------------------------|---|---|---|---|---|---|
| 10 | <b>Inputs</b>              |                                            |   |   |   |   |   |   |
| 11 |                            |                                            |   |   |   |   |   |   |
| 12 | <b>Sampling Parameters</b> |                                            |   |   |   |   |   |   |
| 13 | Directories                | C:\Users\Nicole\Documents\Frontline\TextMi |   |   |   |   |   |   |
| 14 | # Files written            | 300                                        |   |   |   |   |   |   |
| 15 | Write file content?        | FALSE                                      |   |   |   |   |   |   |
| 16 | Sample?                    | TRUE                                       |   |   |   |   |   |   |
| 17 | Sampling method            | Simple random sampling                     |   |   |   |   |   |   |
| 18 | Desired sample size        | 300                                        |   |   |   |   |   |   |
| 19 | Random seed                | 12345                                      |   |   |   |   |   |   |
| 20 |                            |                                            |   |   |   |   |   |   |
| 21 | <b>Text Data</b>           |                                            |   |   |   |   |   |   |
| 22 |                            |                                            |   |   |   |   |   |   |
| 23 | Doc ID                     | TextVar                                    |   |   |   |   |   |   |
| 24 | Doc 1                      | os\101645                                  |   |   |   |   |   |   |
| 25 | Doc 2                      | nics\53662                                 |   |   |   |   |   |   |
| 26 | Doc 3                      | nics\53605                                 |   |   |   |   |   |   |
| 27 | Doc 4                      | nics\53595                                 |   |   |   |   |   |   |
| 28 | Doc 5                      | os\102911                                  |   |   |   |   |   |   |
| 29 | Doc 6                      | os\102901                                  |   |   |   |   |   |   |

Here is an example of a document that appeared in the Electronics newsgroup. Note the appearance of email addresses, “From” and “Subject” lines. All three appear in each document.



The selected file paths are now in random order, but we will need to categorize the “Autos” and “Electronics” files in order to be able to identify them later. To do this, we’ll use Excel to sort the rows by the file path: Select columns B through D and rows 18 through 317, then choose **Sort** from the Data tab. In the Sort dialog, select column d, where the file paths are located, and click OK.





The file paths should now be sorted between *Electronics* and *Autos* files.

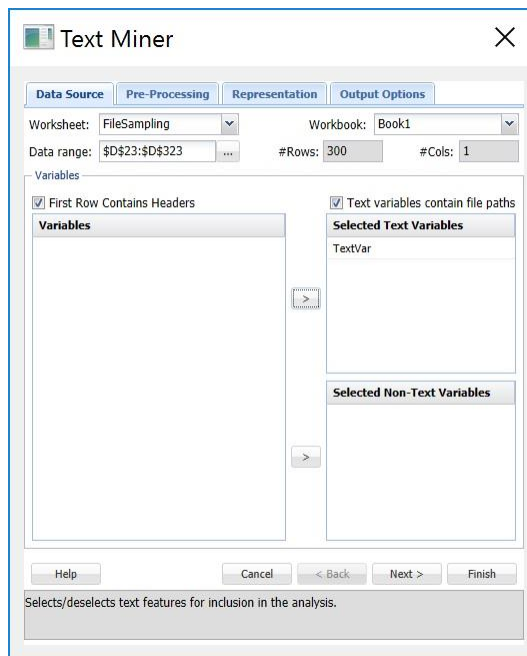
|    | A | B | C       | D                                                                |
|----|---|---|---------|------------------------------------------------------------------|
| 22 |   |   |         |                                                                  |
| 23 |   |   |         |                                                                  |
| 24 |   |   | Doc ID  | TextVar                                                          |
| 25 |   |   | Doc 108 | C:\Users\Nicole\Documents\Frontline\TextMiner Files\Autos\101553 |
| 26 |   |   | Doc 88  | C:\Users\Nicole\Documents\Frontline\TextMiner Files\Autos\101562 |
| 27 |   |   | Doc 125 | C:\Users\Nicole\Documents\Frontline\TextMiner Files\Autos\101564 |
| 28 |   |   | Doc 223 | C:\Users\Nicole\Documents\Frontline\TextMiner Files\Autos\101566 |
| 29 |   |   | Doc 30  | C:\Users\Nicole\Documents\Frontline\TextMiner Files\Autos\101567 |
| 30 |   |   | Doc 194 | C:\Users\Nicole\Documents\Frontline\TextMiner Files\Autos\101568 |

## Using Text Miner

Click the **Text** icon to bring up the *Text Miner* dialog.

### Data Source Tab

Confirm that **XLM\_SampleFiles** is selected for *Worksheet*. Select **TextVar** in the *Variables* list box, and click the upper > button to move it to the *Selected Text Variables* list box. By doing so, we are selecting the text in the documents as input to the Text Miner model. Ensure that “Text variables contain file paths” is **checked**.



Click the **Next** button, or click the **Pre-Processing** tab at the top. b

## Pre-Processing Tab

Leave the default setting for *Analyze all terms* selected under *Mode*. When this option is selected, Text Miner will examine all terms in the document. A “term” is defined as an individual entity in the text, which may or may not be an English word. A term can be a word, number, email, url, etc. terms are separated by all possible delimiting characters (i.e. \, ?, ' , ` , ~ , | , \r , \n , \t , : , ! , @ , # , \$ , % , ^ , & , \* , ( , ) , [ , ] , { , } , < , \_ , ; , = , - , + , \) with some exceptions related to stopwords, synonyms, exclusion terms and boilerplate normalization (URLs, emails, monetary amounts, etc.). Text Miner will not tokenize on these delimiters.

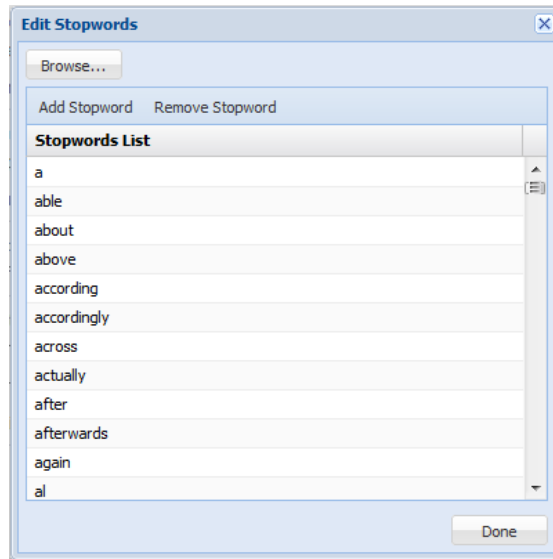
Note: Exceptions are related not to how terms are separated but as to whether they are split based on the delimiter. For example: URL's contain many characters such as "/", ";", etc. Text Miner will not tokenize on these characters in the URL but will consider the URL as a whole and will remove the URL if selected for removal. (See below for more information.)

If *Analyze specified terms only* is selected, the *Edit Terms* button will be enabled. If you click this button, the *Edit Exclusive Terms* dialog opens. Here you can add and remove terms to be considered for text mining. All other terms will be disregarded. For example, if we wanted to mine each document for a specific part name such as “alternator” we would click *Add Term* on the *Edit Exclusive Terms* dialog, then replace “New term” with “alternator” and click *Done* to return to the Pre-Processing dialog. During the text mining process, Text Miner would analyze each document for the term “alternator”, excluding all other terms.

Leave both Start term/phrase and End term/phrase empty under Text Location. If this option is used, text appearing before the first occurrence of the Start Phrase will be disregarded and similarly, text appearing after End Phrase (if used) will be disregarded. For example, if text mining the transcripts from a Live Chat service, you would not be particularly interested in any text appearing before the heading “Chat Transcript” or after the heading “End of Chat Transcript”. Thus you would enter “Chat Transcript” into the Start Phrase field and “End of Chat Transcript” into the End Phrase field.

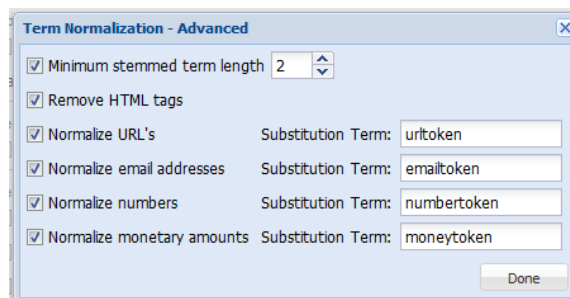
Leave the default setting for *Stopword removal*. Click Edit to view a list of commonly used words that will be removed from the documents during pre-processing. To remove a word from the Stopword list, simply highlight the desired word, then click **Remove Stopword**. To add a new word to the list, click **Add Stopword**, a new term “stopword” will be added. Double click to edit.

Text Miner also allows additional stopwords to be added or existing to be removed via a text document (\*.txt) by using the Browse button to navigate to the file. Terms in the text document can be separated by a space, a comma, or both. If we were supplying our three terms in a text document, rather than in the *Edit Stopwords* dialog, the terms could be listed as: subject emailterm from or subject,emailterm,from or subject, emailterm, from. If we had a large list of additional stopwords, this would be the preferred way to enter the terms.



Click **Advanced** in the *Term Normalization* group to open the *Term Normalization – Advanced* dialog. **Select all options** as shown below. Then click **Done**. This dialog allows us to indicate to Text Miner, that

- If stemming reduced term length to 2 or less characters, disregard the term (*Minimum stemmed term length*).
- HTML tags, and the text enclosed, will be removed entirely. HTML tags and text contained inside these tags often contain technical, computer-generated information that is not typically relevant to the goal of the text mining application.
- URLs will be replaced with the term, “urltoken”. Specific form of URLs do not normally add any meaning, but it is sometimes interesting to know how many URLs are included in a document.
- Email addresses will be replaced with the term, “emailtoken”. Since the documents in our collection all contain a great many email addresses (and the distinction between the different emails often has little use in Text Mining), these email addresses will be replaced with the term “emailtoken”.
- Numbers will be replaced with the term, “numbertoken”.
- Monetary amounts will be substituted with the term, “moneytoken”.



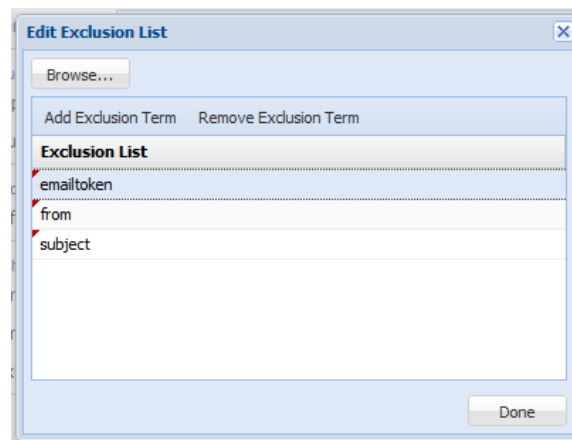
Recall that when we inspected an email from the document collection we saw several terms such as “subject”, “from” and email addresses. Since all of our documents contain these terms, including them in the analysis will not provide any benefit and could bias the analysis. As a result, we will exclude these terms

from all documents by selecting **Exclusion list** then clicking **Edit**. The *Edit Exclusion List* dialog opens. Click **Add Exclusion Term**. The label “exclusionterm” is added. Click to edit and change to “subject”. Then repeat these same steps to add the term “from”.

We can take the email issue one step further and completely remove the term “emailtoken” from the collection. Click **Add Exclusion Term** and edit “exclusionterm” to “emailtoken”.

To remove a term from the exclusion list, highlight the term and click *Remove Exclusion Term*.

We could have also entered these terms into a text document (\*.txt) and added the terms all at once by using the Browse button to navigate to the file and import the list. Terms in the text document can be separated by a space, a comma, or both. If, for example we were supplying excluded terms in a document rather than in the Edit Exclusion List dialog, we would enter the terms as: subject emailtoken from, or subject,emailtoken,from, or subject, emailtoken, from. If we had a large list of terms to be excluded, this would be the preferred way to enter the terms.



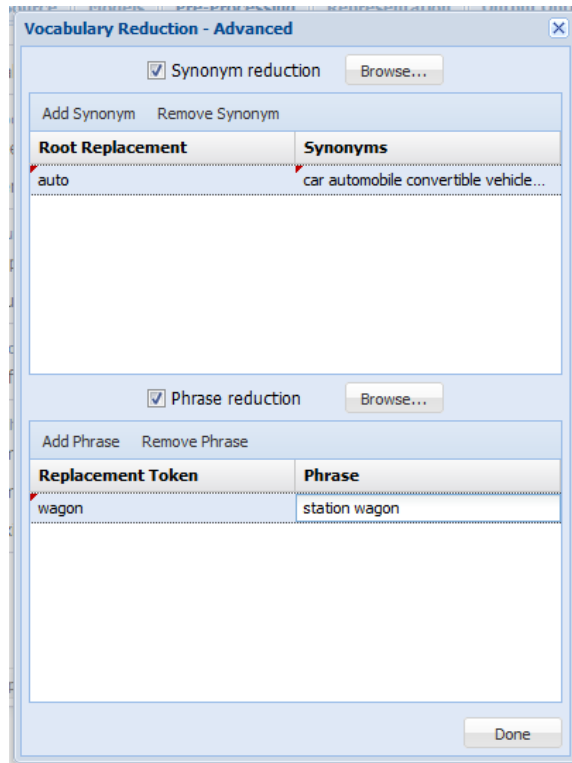
Click **Done** to close the dialog and return to Pre-Processing.

Text Miner also allows the combining of synonyms and full phrases by clicking **Advanced** within *Vocabulary Reduction*. Select **Synonym reduction** at the top of the dialog to replace synonyms such as “car”, “automobile”, “convertible”, “vehicle”, “sedan”, “coupe”, “subcompact”, and “jeep” with “auto”. Click **Add Synonym** and replace “rootterm” with “auto” then replace “synonym list” with “car, automobile, convertible, vehicle, sedan, coupe, subcompact, jeep” (without the quotes). During pre-processing, Text Miner will replace the terms “car”, “automobile”, “convertible”, “vehicle”, “sedan”, “coupe”, “subcompact” and “jeep” with the term “auto”. To remove a synonym from the list, highlight the term and click *Remove Synonym*.

If adding synonyms from a text file, each line must be of the form rootterm:synonymlist or using our example: auto:car automobile convertible vehicle sedan coup or auto:car,automobile,convertible,vehicle,sedan,coup. Note separation between the terms in the synonym list be either a space, a comma or both. If we had a large list of synonyms, this would be the preferred way to enter the terms.

Text Miner also allows the combining of words into phrases that indicate a singular meaning such as “station wagon” which refers to a specific type of car rather than two distinct tokens – station and wagon. To add a phrase in the *Vocabulary Reduction – Advanced* dialog, select **Phrase reduction** and click

**Add Phrase.** The term “phrasetoken” will be appear, click to edit and enter “wagon”. Click “phrase” to edit and enter “station wagon”. If supplying phrases through a text file (\*.txt), each line of the file must be of the form phrasetoken:phrase or using our example, wagon:station wagon. If we had a large list of phrases, this would be the preferred way to enter the terms.



Enter **200** for *Maximum Vocabulary Size*. Text Miner will reduce the number of terms in the final vocabulary to the top 200 most frequently occurring in the collection.

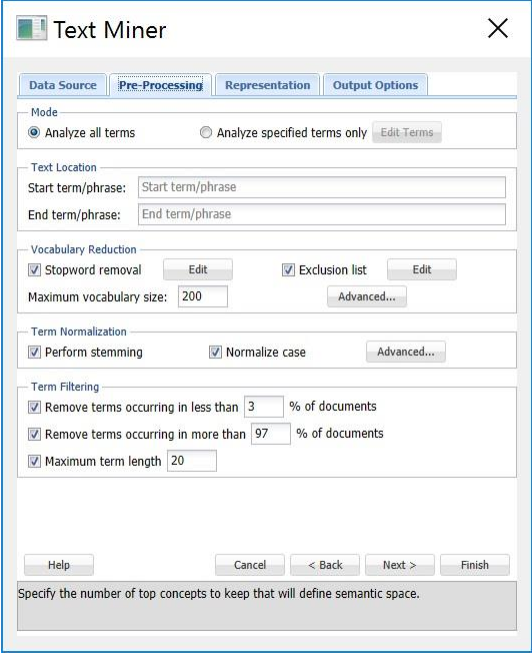
Leave *Perform stemming* at the selected default. Stemming is the practice of stripping words down to their “stems” or “roots”, for example, stemming terms such as “argue”, “argued”, “argues”, “arguing”, and “argus” would result in the stem “argu. However “argument” and “arguments” would stem to “argument”. The stemming algorithm utilized in Text Miner is “smart” in the sense that while “running” would be stemmed to “run”, “runner” would not. . Text Miner uses the Porter Stemmer 2 algorithm for the English Language. For more information on this algorithm, please see the Webpage: <http://tartarus.org/martin/PorterStemmer/>

Leave the default selection for *Normalize case*. When this option is checked, Text Miner converts all text to a consistent (lower) case, so that Term, term, TERM, etc. are all normalized to a single token “term” before any processing, rather than creating three independent tokens with different case. This simple method can dramatically affect the frequency distributions of the corpus, leading to biased results.

Enter **3** for *Remove terms occurring in less than \_% of documents* and **97** for *Remove terms occurring in more than \_% of documents*. For many text mining applications, the goal is to identify terms that are useful for discriminating between documents. If a particular term occurs in all or almost all documents, it may not be possible to highlight the differences. If a term occurs in very few

documents, it will often indicate great specificity of this term, which is not very useful for some Text Mining purposes.

Enter **20** for *Maximum term length*. Terms that contain more than 20 characters will be excluded from the text mining analysis and will not be present in the final reports. This option can be extremely useful for removing some parts of text which are not actual English words, for example, URLs or computer-generated tokens, or to exclude very rare terms such as Latin species or disease names, i.e. Pneumonoultramicroscopicsilicovolcanoconiosis.



The screenshot shows the 'Text Miner' application window with the 'Pre-Processing' tab selected. The 'Mode' section has 'Analyze all terms' selected. The 'Text Location' section has empty input fields for 'Start term/phrase' and 'End term/phrase'. The 'Vocabulary Reduction' section has 'Stopword removal' and 'Exclusion list' checked, with a 'Maximum vocabulary size' of 200. The 'Term Normalization' section has 'Perform stemming' and 'Normalize case' checked. The 'Term Filtering' section has 'Remove terms occurring in less than 3 % of documents', 'Remove terms occurring in more than 97 % of documents', and 'Maximum term length 20' checked. At the bottom, there are buttons for 'Help', 'Cancel', '< Back', 'Next >', and 'Finish'. A note at the very bottom says 'Specify the number of top concepts to keep that will define semantic space.'

Click **Next** to advance to *the Representation* tab or simply click **Representation** at the top.

## **Representation Tab**

Keep the default selection of *TF-IDF* (Term Frequency – Inverse Document Frequency) for *Term-Document Matrix Scheme*. A term-document matrix is a matrix that displays the frequency-based information of terms occurring in a document or collection of documents. Each column is assigned a term and each row a document. If a term appears in a document, a weight is placed in the corresponding column indicating the term's importance or contribution. Text Miner offers four different commonly used methods of weighting scheme used to represent each value in the matrix: presence/Absence, Term Frequency, TF-IDF (the default) and Scaled term frequency. If *Presence/Absence* is selected, Text Miner will enter a 1 in the corresponding row/column if the term appears in the document and 0 otherwise. This matrix scheme does not take into account the number of times the term occurs in each document. If *Term Frequency* is selected, Text Miner will count the number of times the term appears in the document and enter this value into the corresponding row/column in the matrix. The default setting – Term Frequency – Inverse Document Frequency (*TF-IDF*) is the product of scaled term frequency and inverse document frequency. Inverse document frequency is calculated by taking the logarithm of the total number of documents divided by the number of documents that contain the term. A high value for TF-IDF indicates that a term that does not occur frequently in the collection of documents taken as a whole, appears quite

frequently in the specified document. A TF-IDF value close to 0 indicates that the term appears frequently in the collection or rarely in a specific document. If *Scaled term frequency* is selected, Text Miner will normalize (bring to the same scale) the number of occurrences of a term in the documents (see the table below).

It's also possible to create your own scheme by clicking the *Advanced* command button to open the *Term Document Matrix – Advanced* dialog. Here users can select their own choices for local weighting, global weighting, and normalization. Please see the table below for definitions regarding options for Term Frequency, Document Frequency and Normalization.

| Local Weighting |                                                                                                  | Global Weighting |                                                       | Normalization |                                   |
|-----------------|--------------------------------------------------------------------------------------------------|------------------|-------------------------------------------------------|---------------|-----------------------------------|
| Binary          | $lw_{td} = \begin{cases} 1, & \text{if } tf_{td} > 0 \\ 0, & \text{if } tf_{td} = 0 \end{cases}$ | None             | $gw_t = 1$                                            | None          | $n_d = 1$                         |
| Raw Frequency   | $lw_{td} = tf_{td}$                                                                              | Inverse          | $gw_t = \log_2 \frac{N}{1 + df_t}$                    | Cosine        | $n_d = \frac{1}{\ \bar{g}_d\ _2}$ |
| Logarithmic     | $lw_{td} = \log(1 + tf_{td})$                                                                    | Normal           | $gw_t = \frac{1}{\sqrt{\sum_d tf_{td}^2}}$            |               |                                   |
| Augnorm         | $lw_{td} = \frac{\left(\frac{tf_{td}}{\max_t tf_{td}}\right) + 1}{2}$                            | GF-IDF           | $gw_t = \frac{cf_t}{df_t}$                            |               |                                   |
|                 |                                                                                                  | Entropy          | $gw_t = 1 + \sum_d \frac{p_{td} \log p_{td}}{\log N}$ |               |                                   |
|                 |                                                                                                  | IDF probability  | $gw_t = \log_2 \frac{N}{1 + df_t}$                    |               |                                   |

Notations:

- $tf_{td}$  – frequency of term  $t$  in a document  $d$ ;
- $df_t$  – document frequency of term  $t$ ;
- $lw_{td}$  – local weighting of term  $t$  in a document  $d$ ;
- $gw_{td}$  – global weighting of term  $t$  in a document  $d$ ;
- $n_d$  – normalization of vector of terms representing the document  $d$ ;
- $N$  – total number of documents in the collection;
- $cf_t$  – collection frequency of term  $t$ ;
- $p_{td}$  – estimated probability of term  $t$  to appear in a document  $d$  ( $p_{td} = tf_{td}/cf_t$ );
- $\bar{g}_d$  – vector of terms representing the document  $d$ .

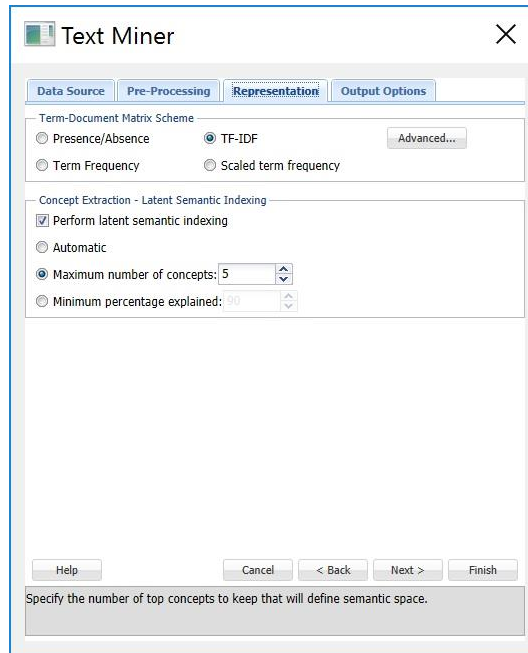
Finally, the element  $T_{td}$  of Term-Document Matrix is computed as  $T_{td} = lw_{td} * gw_t * n_d, \forall t, d$

Leave *Perform latent semantic indexing* selected (the default). When this option is selected, Text Miner will use Latent Semantic Indexing (LSI) to detect patterns in the associations between terms and concepts to discover the meaning of the document.

The statistics produced and displayed in the Term-Document Matrix contain basic information on the frequency of terms appearing in the document collection. With this information we can “rank” the significance or importance of these terms relative to the collection and particular document. Latent Semantic Indexing, in comparison, uses singular value decomposition (SVD) to map the terms and documents into a common space to find patterns and relationships. For example: if we inspected our document collection, we might find that each time the term “alternator” appeared in an automobile document, the document also included the terms “battery” and “headlights”. Or each time the term “brake” appeared in an automobile document, the terms “pads” and “squeaky” also appeared. However there is no detectable pattern regarding the use of the terms “alternator” and “brake”. Documents including “alternator” might not include “brake” and documents including “brake” might not include “alternator”. Our four terms, battery, headlights, pads, and squeaky describe two different automobile repair issues: failing brakes and a bad alternator. Latent Semantic Indexing will attempt to 1. Distinguish between these two different topics, 2. Identify the documents that deal with faulty brakes, alternator problems or both and 3. Map the terms into a common semantic space using singular value decomposition. SVD is a tool used by Text Miner to extract concepts that explain the main dimensions of meaning of the documents in the collection. The results of LSA are usually hard to examine because the construction of the concept representations will not be fully explained. Interpreting these results is actually more of an art, than a science. However, Text Miner provides several visualizations that simplify this process greatly.

Select **Maximum number of concepts** and leave the default setting of 5. Doing so will tell Text Miner to retain the top 5 most significant concepts. If *Automatic* is selected, Text Miner will calculate the importance of each concept, take the difference between each and report any concepts above the largest difference. For example if three concepts were identified (Concept1, Concept2, and Concept3) and given importance factors of 10, 8, and 2, respectively, Text Miner would keep Concept1 and Concept2 since the difference between Concept2 and Concept 3 (8-2=6) is larger than the difference between Concept1 and Concept2 (10-8=2). If *Minimum percentage explained* is selected, Text Miner will identify the concepts with singular values that, when taken together, sum to the minimum percentage explained, 90% is the default.





Click **Next** or the **Output Options** tab.

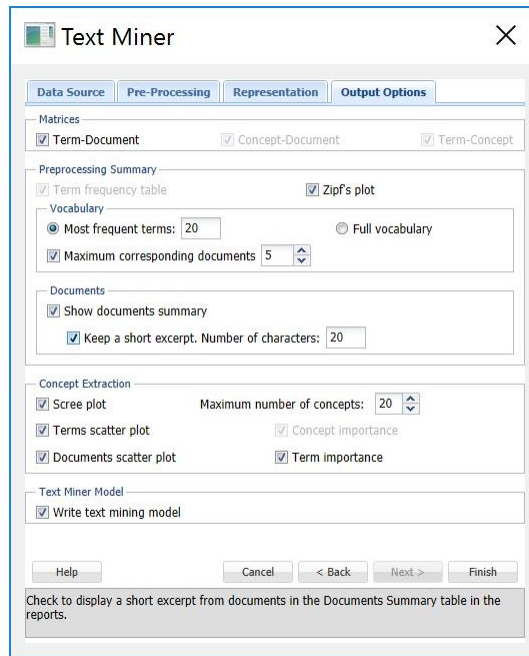
## Options Tab

Keep *Term-Document* and *Concept-Document* selected under *Matrices* (the default) and select *Term-Concept* to print each matrix in the output. The *Term-Document* matrix displays the terms across the top of the matrix and the documents down the left side of the matrix. The *Concept – Document* and *Term – Concept* matrices are output from the *Perform latent semantic indexing* option that we selected on the *Representation* tab. In the first matrix, *Concept – Document*, 20 concepts will be listed across the top of the matrix and the documents will be listed down the left side of the matrix. The values in this matrix represent concept coordinates in the identified semantic space. In the *Term-Concept* matrix, the terms will be listed across the top of the matrix and the concepts will be listed down the left side of the matrix. The values in this matrix represent terms in the extracted semantic space.

Keep *Term frequency table* selected (the default) under *Preprocessing Summary* and select *Zipf's plot*. Increase the *Most frequent terms* to 20 and select *Maximum corresponding* documents. The Term frequency table will include the top 20 most frequently occurring terms. The first column, *Collection Frequency*, displays the number of times the term appears in the collection. The 2<sup>nd</sup> column, *Document Frequency*, will display the number of documents that include the term. The third column, *Top Documents*, will display the top 5 documents where the corresponding term appears the most frequently. The *Zipf Plot* graphs the document frequency against the term ranks in descending order of frequency.. Zipf's law states that the frequency of terms used in a free-form text drops exponentially, i.e. that people tend to use a relatively small number of words extremely frequently and use a large number of words very rarely.

Keep *Show documents summary* selected and check *Keep a short excerpt*. under *Documents*. Text Miner will produce a table displaying the document ID, length of the document, number of terms and 20 characters of the text of the document.

Select **all plots** under *Concept Extraction* to produce various plots in the output. Select *Write text mining model* under *Text Miner Model* to write the model to an output sheet.



Click the **Finish** button to run the Text Mining analysis. Result worksheets are inserted to the right.

## Output Results

The Output Navigator appears at the top of each output worksheet. Clicking any of these links will allow you to "jump" to the desired output.

|   | B                         | C               | D                       | E                   | F      | G | H | I | J | K |
|---|---------------------------|-----------------|-------------------------|---------------------|--------|---|---|---|---|---|
| 1 | Data Science: Text Mining |                 |                         |                     |        |   |   |   |   |   |
| 2 |                           |                 |                         |                     |        |   |   |   |   |   |
| 3 | <b>Output Navigator</b>   |                 |                         |                     |        |   |   |   |   |   |
| 4 | Concept Importance        | Term Importance | Concept-Document Matrix | Term-Concept Matrix | Inputs |   |   |   |   |   |
| 5 | Term Count Info           | Document Info   | Term-Document Matrix    | Top Terms Info      |        |   |   |   |   |   |
| 6 |                           |                 |                         |                     |        |   |   |   |   |   |

### Term Count and Document Info

Select the *TM\_Output* tab. The *Term Count* table shows that the original term count in the documents was reduced by 14.26% by the removal of stopwords, excluded terms, synonyms, phrase removal and other specified preprocessing procedures.

|    | B                      | C                | D             | E            | F          | G |
|----|------------------------|------------------|---------------|--------------|------------|---|
| 44 | <b>Term Count Info</b> |                  |               |              |            |   |
| 45 |                        |                  |               |              |            |   |
| 46 | Text Var               | Original (Total) | Final (Total) | Reduction, % | Vocabulary |   |
| 47 | TextVar                | 70803            | 10099         | 14.26351991  | 200        |   |
| 48 |                        |                  |               |              |            |   |

Scroll down to the *Documents* table. This table lists each Document with its length, number of terms, and if *Keep a short excerpt* is selected on the *Output Options* tab and a value is present for *Number of characters*, then an excerpt from each document will be displayed.

49 **Document Info**

| Document ID | # Characters | # Terms | Excerpt: TextVar          |
|-------------|--------------|---------|---------------------------|
| 101553      |              | 287     | 50 ; netops@tekgen....    |
| 101562      |              | 535     | 91 i: stlucas@gdwest...   |
| 101564      |              | 769     | 140 edwards@world....     |
| 101566      |              | 1125    | 195 : tbigham@shears...   |
| 101567      |              | 981     | 160 om: silver@xrtll.u... |
| 101568      |              | 1146    | 177 om: silver@xrtll.u... |

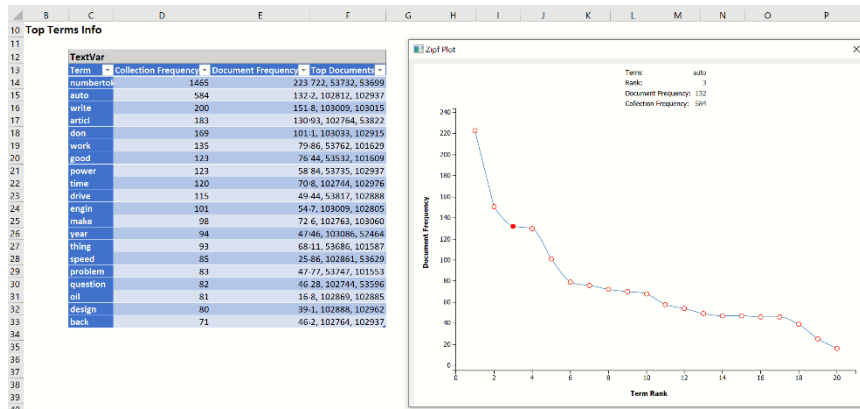
### Term-Document Matrix

Click *TM\_TDM*, to display the Term – Document Matrix. As discussed above, this matrix lists the 200 most frequently appearing terms across the top and the document IDs down the left. A portion of this table is shown below. If a term appears in a document, a weight is placed in the corresponding column indicating the importance of the term using our selection of TF-IDF on the Representation dialog.

| Doc ID | ago      | air | amp | appreci | area | articl   | auto     | avail | away | back | bad | base | believ | best | better   | big | bit | book | box |
|--------|----------|-----|-----|---------|------|----------|----------|-------|------|------|-----|------|--------|------|----------|-----|-----|------|-----|
| 101553 | 0        | 0   | 0   | 0       | 0    | 0        | 0        | 0     | 0    | 0    | 0   | 0    | 0      | 0    | 0        | 0   | 0   | 0    | 0   |
| 101562 | 0        | 0   | 0   | 0       | 0    | 0        | 0.563829 | 0     | 0    | 0    | 0   | 0    | 0      | 0    | 0        | 0   | 0   | 0    | 0   |
| 101564 | 0        | 0   | 0   | 0       | 0    | 0.574331 | 0        | 0     | 0    | 0    | 0   | 0    | 0      | 0    | 0        | 0   | 0   | 0    | 0   |
| 101566 | 0        | 0   | 0   | 0       | 0    | 0        | 1.30917  | 0     | 0    | 0    | 0   | 0    | 0      | 0    | 1.780202 | 0   | 0   | 0    | 0   |
| 101567 | 1.877077 | 0   | 0   | 0       | 0    | 0        | 0.893648 | 0     | 0    | 0    | 0   | 0    | 0      | 0    | 0        | 0   | 0   | 0    | 0   |
| 101568 | 0        | 0   | 0   | 0       | 0    | 0        | 0        | 0     | 0    | 0    | 0   | 0    | 0      | 0    | 0        | 0   | 0   | 0    | 0   |

### Vocabulary Matrix

Click *TM\_Vocabulary* to view the Final List of Terms table. This table contains the top 20 terms occurring in the document collection, the number of documents that include the term and the top 5 document IDs where the corresponding term appears most frequently. In this list we see terms such as “car”, “power”, “engine”, “drive”, and “dealer” which suggests that many of the documents, even the documents from the electronic newsgroup, were related to autos.



When you click on the *TM\_Vocabulary* tab, the Zipf Plot opens. We see that our collection of documents obey the power law stated by Zipf (see above). As we move from left to right on the graph, the documents that contain the most frequently appearing terms (when ranked from most frequent to least frequent) drop quite steeply. Hover over each data point to see the detailed information about the term corresponding to this data point.

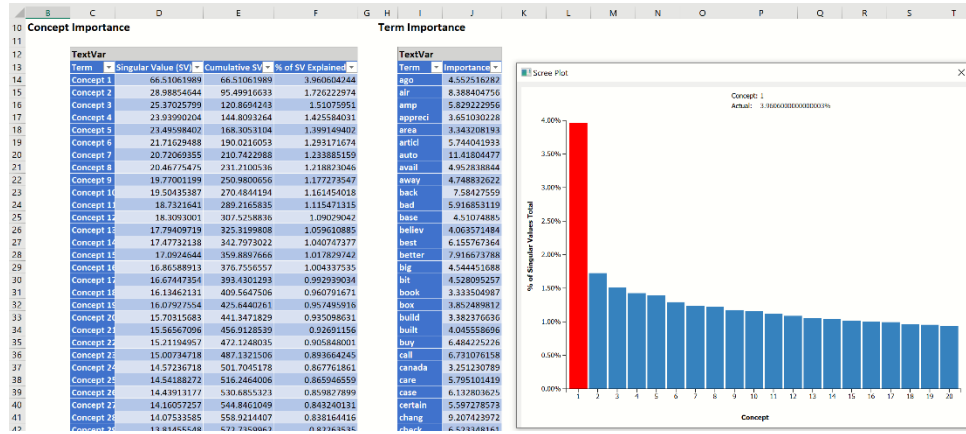
Note: To view charts in the Data Science Cloud app, click *Charts* on the Ribbon, select the desired worksheet, in this case *TM\_Vocabulary*, then select the desired chart.

The term “numbertoken” is the most frequently occurring term in the document collection appearing in 223 documents (out of 300), 1,083 times total. Compare

this to a less frequently occurring term such as "thing" which appears in only 64 documents and only 82 times total.

### Concept Importance

Click *TM\_LSASummary* to view the Concept Importance and Term Importance tables. The first table, the Concept Importance table, lists each concept, its singular value, the cumulative singular value and the % singular value explained. (The number of concepts extracted is the minimum of the number of documents (985) and the number of terms (limited to 200).) These values are used to determine which concepts should be used in the Concept – Document Matrix, Concept – Term Matrix and the Scree Plot according to the Users selection on the Representation tab. In this example, we entered “20” for *Maximum number of concepts*.



The Term Importance table lists the 200 most important terms. (To increase the number of terms from 200, enter a larger value for Maximum Vocabulary on the Pre-processing tab of Text Miner.)

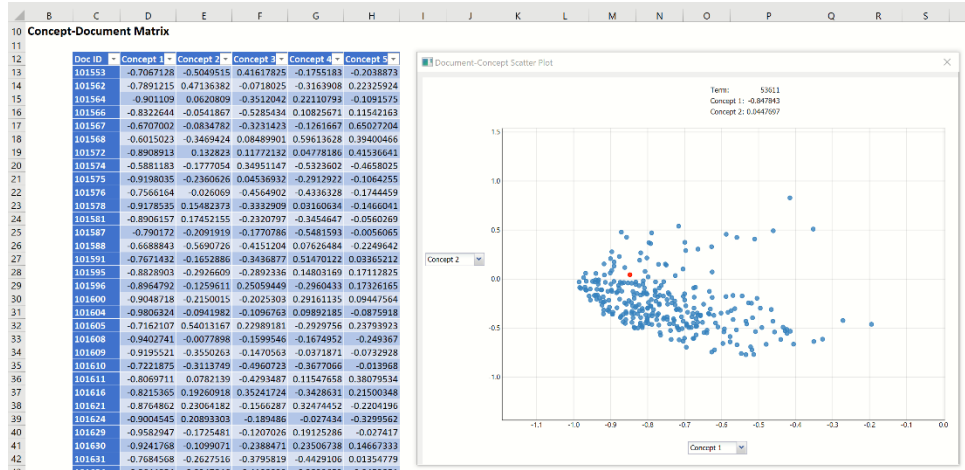
When you click the *TM\_LSASummary* tab, the Scree Plot opens. This plot gives a graphical representation of the contribution or importance of each concept. The largest “drop” or “elbow” in the plot appears between the 1<sup>st</sup> and 2<sup>nd</sup> concept. This suggests that the first top concept explains the leading topic in our collection of documents. Any remaining concepts have significantly reduced importance. However, we can always select more than 1 concept to increase the accuracy of the analysis – it is advised to examine the Concept Importance table and the “Cumulative Singular Value” in particular to identify how many top concepts capture enough information for your application.

### Concept Document Matrix

Click *TM\_LSA\_CDM* to display the Concept – Document Matrix. This matrix displays the top concepts (as selected on the Representation tab) along the top of the matrix and the documents down the left side of the matrix.

When you click on the *TM\_LSA\_CDM* tab, the Concept-Document Scatter Plot opens. This graph is a visual representation of the Concept – Document matrix. Note that Analytic Solver Data Science normalizes each document representation so it lies on a unit hypersphere. Documents that appear in the middle of the plot, with concept coordinates near 0 are not explained well by either of the shown concepts. The further the magnitude of coordinate from zero, the more effect that particular concept has for the corresponding document. In fact, two documents placed at extremes of a concept (one close to -1 and other to +1) indicates strong differentiation between these documents in terms of

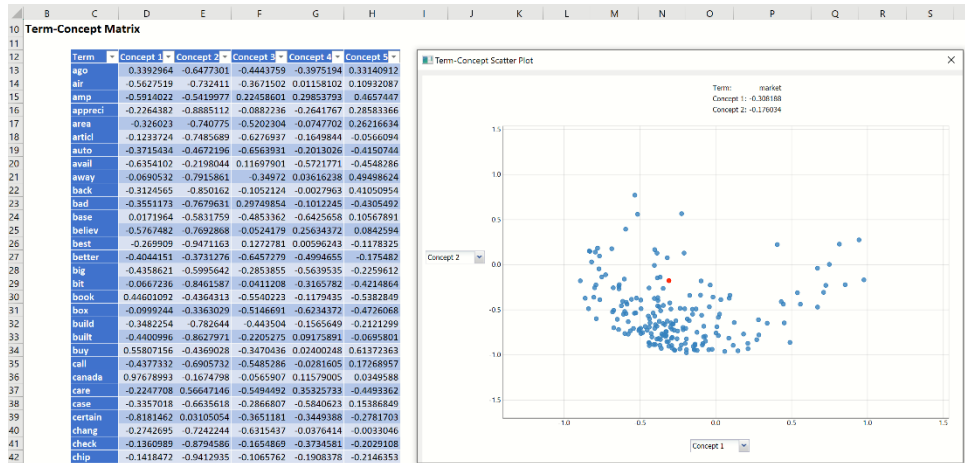
the extracted concept. This provides means for understanding actual meaning of the concept and investigating which concepts have the largest discriminative power, when used to represent the documents from the text collection.



You can examine all extracted concepts by changing the axes on a scatter plot - click the down pointing arrow next to Concept 1 or the concept on the Y axis by clicking the right pointing arrow next to Concept 2. Use your touchscreen or your mouse scroll wheel to zoom in and out.

### Term-Concept Matrix

Double click *TM\_LSA\_CTM* to display the Concept – Term Matrix which lists the top 5 most important concepts along the top of the matrix and the top 200 most frequently appearing terms down the side of the matrix.



When you click on the *TM\_LSA-CTM* tab, the Term-Concept Scatter Plot opens. This graph is a visual representation of the Concept – Term Matrix. It displays all terms from the final vocabulary in terms of two concepts. Similarly to the Concept-Document scatter plot, the Concept-Term scatter plot visualizes the distribution of vocabulary terms in the semantic space of meaning extracted with LSA. The coordinates are also normalized, so the range of axes is always [-1, 1], where extreme values (close to +/-1) highlight the importance or “load” of each term to a particular concept. The terms appearing in a zero-neighborhood of concept range do not contribute much to a concept definition. In our example, if we identify a concept having a set of terms that can be divided into two groups: one related to “Autos” and other to “Electronics”, and these groups are

distant from each other on the axis corresponding to this concept, this would definitely provide an evidence that this particular concept “caught” some pattern in the text collection that is capable of discriminating the topic of article. Therefore, Term-Concept scatter plot is an extremely valuable tool for examining and understanding the main topics in the collection of documents, finding similar words that indicate similar concept, or the terms explaining the concept from “opposite sides” (e.g. *term1* can be related to cheap affordable electronics and *term2* can be related to expensive luxury electronics)

Recall that if you want to examine different pair of concepts, click the down pointing arrow next to Concept 1 and the right pointing arrow next to Concept 2 to change the concepts on either axis. Use your touchscreen or mouse wheel to scroll in or out.

### **Stored PMML models for TFIDF and LSA**

Since "Write text mining model" on the Output Options tab, two more tabs are created containing PMML models for the TFIDF and LSA models. These models can be used when scoring a series of new documents. For more information on how to process this new data using these two saved models, see the Text Mining chapter within the Data Science Reference Guide.

## **Classification with Concept Document Matrix**

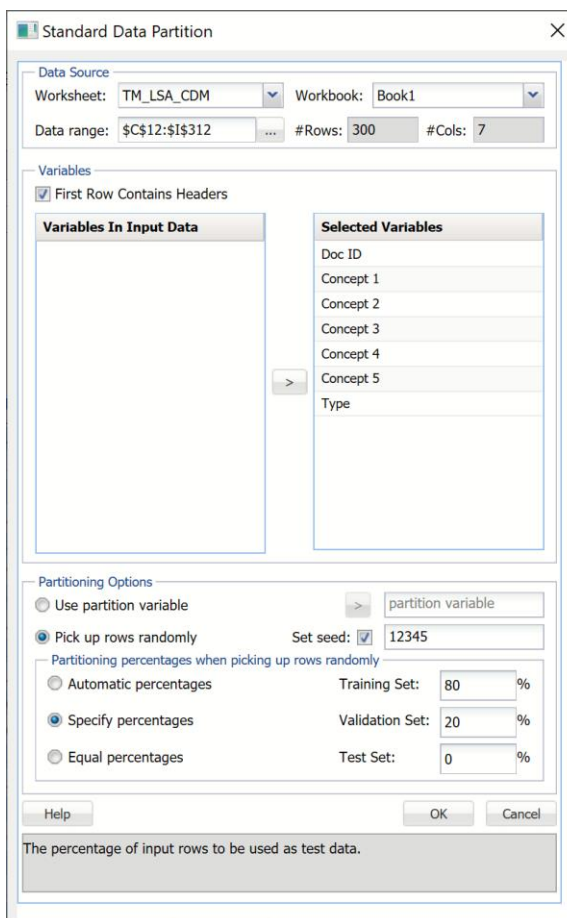
From here, we can use any of the six classification algorithms to classify our documents according to some term or concept using the Term – Document matrix, Concept – Document matrix or Concept – Term matrix where each document becomes a “record” and each concept becomes a “variable”. If wanting to classify documents based on a binary variable such as Auto email/non-Auto email, then we would use either the Term – Document or Concept – Document matrix. If wanting to cluster terms or classify terms, then we would use the Term-Concept matrix. We could even use the transpose of the Term – Document matrix where each term would become a “record” and each column would become a “feature”. See the Analytic Solver Data Science User Guide for an example model that uses the Logistic Regression Classification method to create a classification model using the Concept Document matrix within *TM\_LSA\_CDM*.

In this example, we will use the Logistic Regression Classification method to create a classification model using the Concept Document matrix on *TM\_LSA\_CDM*. Recall that this matrix includes the top twenty concepts extracted from the document collection across the top of the matrix and each document in the sample down the left. Each concept will now become a “feature” and each document will now become a “record”.

First, we’ll need to append a new column with the class that the document is currently assigned: electronics or autos. Since we sorted our documents at the beginning of the example starting with Autos, we can simply enter “Autos” into column I for Document IDs 101553 through 103096 (or cells I13:I162) and enter “Electronics” into column I for Document IDs 52434 through 53879 (or cells I163:I312). Give the column the title "Type".

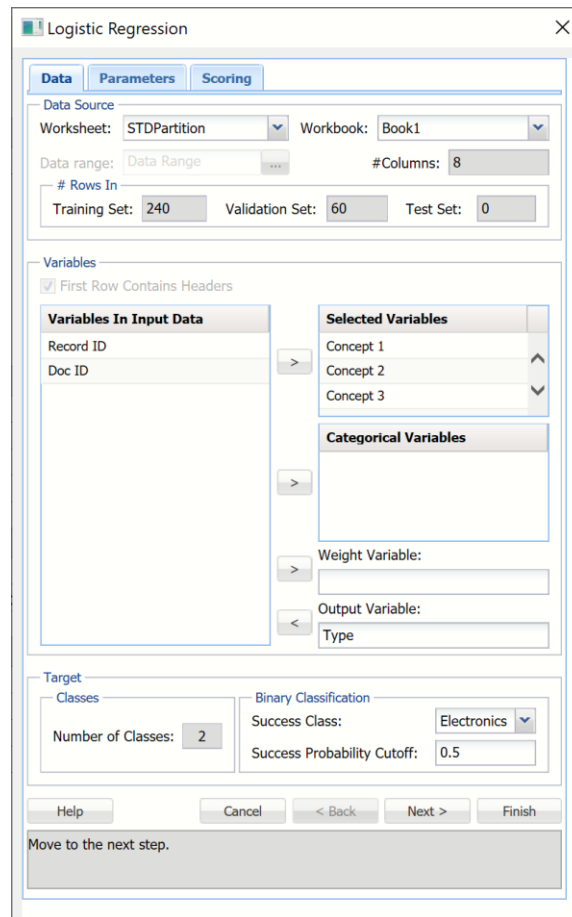
|    | B | C      | D          | E          | F          | G          | H          | I    |
|----|---|--------|------------|------------|------------|------------|------------|------|
| 11 |   |        |            |            |            |            |            |      |
| 12 |   | Doc ID | Concept 1  | Concept 2  | Concept 3  | Concept 4  | Concept 5  | Type |
| 13 |   | 101553 | -0.6445755 | 0.57575418 | -0.3599934 | -0.0637993 | 0.34549084 | Auto |
| 14 |   | 101562 | -0.9218592 | -0.1203888 | 0.36765001 | 0.01563106 | -0.0164715 | Auto |
| 15 |   | 101564 | -0.9369048 | -0.2671295 | 0.21632776 | -0.0217637 | 0.05983167 | Auto |
| 16 |   | 101566 | -0.8098641 | -0.1202806 | 0.4767241  | -0.3015155 | 0.10712239 | Auto |
| 17 |   | 101567 | -0.5862534 | 0.28093583 | 0.46033063 | -0.540465  | -0.2708787 | Auto |
| 18 |   | 101568 | -0.6232086 | 0.16053979 | -0.2514671 | -0.6001768 | 0.40297659 | Auto |
| 19 |   | 101572 | -0.899317  | -0.254814  | -0.2378974 | -0.1653817 | -0.2057972 | Auto |
| 20 |   | 101574 | -0.5659841 | 0.34030445 | 0.02131019 | 0.69445309 | 0.28484325 | Auto |
| 21 |   | 101575 | -0.9166982 | 0.19430103 | 0.03506041 | 0.14282556 | 0.31667504 | Auto |
| 22 |   | 101576 | -0.6895418 | -0.0130039 | 0.56113989 | 0.2137582  | 0.40471278 | Auto |
| 23 |   | 101578 | -0.8016409 | 0.0298667  | 0.28741838 | 0.08933483 | 0.51564502 | Auto |
| 24 |   | 101581 | -0.8629759 | -0.2248615 | 0.37983875 | 0.23915115 | -0.0569139 | Auto |
| 25 |   | 101587 | -0.8411039 | 0.17412812 | 0.48289173 | 0.17033887 | -0.0048915 | Auto |
| 26 |   | 101588 | -0.5702969 | 0.20555356 | 0.45655105 | -0.4723217 | 0.44831072 | Auto |
| 27 |   | 101591 | -0.3174672 | 0.45045962 | -0.0059071 | -0.5249706 | 0.64859203 | Auto |
| 28 |   | 101595 | -0.8283696 | 0.1510668  | 0.34161134 | -0.409604  | -0.0806774 | Auto |
| 29 |   | 101596 | -0.8565701 | 0.42082139 | 0.06176657 | 0.29101869 | -0.02627   | Auto |
| 30 |   | 101600 | -0.8235769 | -0.1874607 | 0.05335228 | -0.4981872 | -0.1885277 | Auto |
| 31 |   | 101604 | -0.9685009 | -0.0157072 | 0.20443555 | -0.1228535 | -0.0698024 | Auto |

First, we'll need to partition our data into two datasets, a training dataset where the model will be "trained" and a validation dataset where the newly created model can be tested, or validated. When the model is being trained, the actual class label assignments are "shown" to the algorithm in order for it to "learn" which variables (or concepts) result in an "auto" or "electronic" assignment. When the model is being validated or tested, the known classification is only used to evaluate the performance of the algorithm. Click **Partition – Standard Partition** on the Text Miner ribbon to open the *Standard Data Partition* dialog. Select all variables in the *Variables In Input Data* list box, then click > to move all to the *Selected Variables* list box. Select **Specify percentages** Under *Partitioning percentages when picking up rows randomly* (at the bottom) and enter **80** for *Training Set*. Automatically, **20** will be entered for *Validation Set*.



Click **Finish** to partition the data into two randomly selected datasets: The Training dataset containing 80% of the “records” (or documents) and the Validation dataset containing 20% of the “records”. (For more information on partitioning, please see the *Standard Partitioning* chapter that appears in the Analytic Solver Data Science Reference Guide.)

Now click Classify – Logistic Regression to open the *Logistic Regression – Step 1 of 3* dialog. Select all 5 concepts under *Variables In Input Data* list box and click > to move them to the *Selected Variables* list box. Doing so selects these variables as inputs to the classification method. Select **Type**, then click the > next to *Output Variable* to add this variable as the Output Variable.



Click **Finish** to accept all defaults and run Logistic Regression.

Select *DA\_ValidationScore* tab and scroll down to the Validation Classification Summary, shown below.

Text Miner used the training dataset to “train” the Logistic Regression model to classify each “record” (or document) as an “autos” or “electronics” document. Afterwards, Text Miner tested the newly created Logistic Regression model on the records in the validation dataset and assigned each record (or document) a classification.



|    | B                                         | C                       | D           | E           | F        |
|----|-------------------------------------------|-------------------------|-------------|-------------|----------|
| 10 | <b>Validation: Classification Summary</b> |                         |             |             |          |
| 11 |                                           |                         |             |             |          |
| 12 |                                           | <b>Confusion Matrix</b> |             |             |          |
| 13 |                                           | Actual\Predicted        | Auto        | Electronics |          |
| 14 |                                           | Auto                    | 23          | 4           |          |
| 15 |                                           | Electronics             | 6           | 27          |          |
| 16 |                                           |                         |             |             |          |
| 17 |                                           | <b>Error Report</b>     |             |             |          |
| 18 |                                           | Class                   | # Cases     | # Errors    | % Error  |
| 19 |                                           | Auto                    | 27          | 4           | 14.81481 |
| 20 |                                           | Electronics             | 33          | 6           | 18.18182 |
| 21 |                                           | Overall                 | 60          | 10          | 16.66667 |
| 22 |                                           |                         |             |             |          |
| 23 |                                           | <b>Metrics</b>          |             |             |          |
| 24 |                                           | Metric                  | Value       |             |          |
| 25 |                                           | Accuracy (#correct)     | 50          |             |          |
| 26 |                                           | Accuracy (%correct)     | 83.33333    |             |          |
| 27 |                                           | Specificity             | 0.851852    |             |          |
| 28 |                                           | Sensitivity (Recall)    | 0.818182    |             |          |
| 29 |                                           | Precision               | 0.870968    |             |          |
| 30 |                                           | F1 score                | 0.84375     |             |          |
| 31 |                                           | Success Class           | Electronics |             |          |
| 32 |                                           | Success Probability     | 0.5         |             |          |
| 33 |                                           |                         |             |             |          |

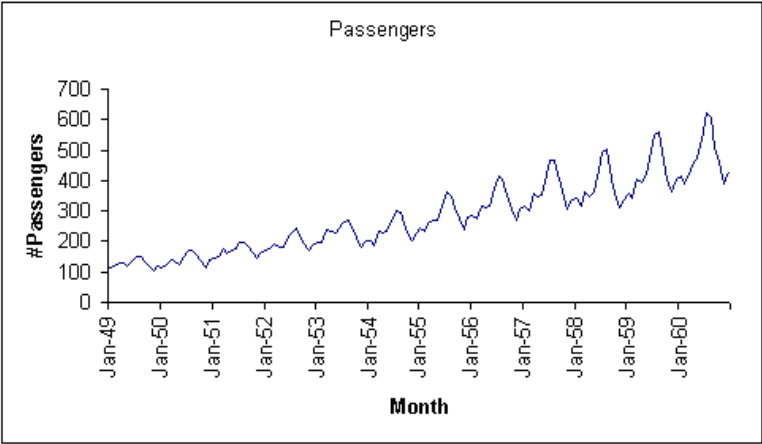
As you can see in the reports above, Logistic Regression was able to correctly classify 50 out of a total of 60 documents in the validation partition, which translates to an overall error of 16.67%, (For more information on how to read the summary report, see the Logistic Regression chapter later on in this guide.)

This concludes our example on how to use Analytic Solver Data Science's Text Miner feature. This example has illustrated how Analytic Solver Data Science provides powerful tools for importing a collection of documents for comprehensive text preprocessing, quantitation, and concept extraction, in order to create a model that can be used to process new documents – all performed without any manual intervention. When using Text Miner in conjunction with our classification algorithms, Analytic Solver Data Science can be used to classify customer reviews as satisfied/not satisfied, distinguish between which products garnered the least negative reviews, extract the topics of articles, cluster the documents/terms, etc. The applications for Text Miner are endless!

# Exploring a Time Series Dataset

## Introduction

Time series datasets contain a set of observations generated sequentially in time. Organizations of all types and sizes utilize time series datasets for analysis and forecasting for predicting next year’s sales figures, raw material demand, monthly airline bookings, etc. .



Example of a time series dataset: Monthly airline bookings.

A time series model is first used to obtain an understanding of the underlying forces and structure that produced the data and then secondly, to fit a model that will predict future behavior. In the first step, the analysis of the data, a model is created to uncover seasonal patterns or trends in the data, for example bathing suit sales in June. In the second step, forecasting, the model is used to predict the value of the data in the future, for example, next year’s bathing suit sales. Separate modeling methods are required to create each type of model.

Analytic Solver Data Science features two techniques for exploring trends in a dataset, ACF (Autocorrelation function) and PACF (Partial autocorrelation function). These techniques help the user to explore various patterns in the data which can be used in the creation of the model. After the data is analyzed, a model can be fit to the data using Analytic Solver Data Science’s ARIMA method.

## Autocorrelation (ACF)

**Autocorrelation (ACF)** is the correlation between neighboring observations in a time series. When determining if an autocorrelation exists, the original time series is compared to the “lagged” series. This lagged series is simply the original series moved one time period forward ( $x_n$  vs  $x_{n+1}$ ). Suppose there are 5 time based observations: 10, 20, 30, 40, and 50. When lag = 1, the original series is moved forward one time period. When lag = 2, the original series is moved forward two time periods.

| Day | Observed Value | Lag-1 | Lag-2 |
|-----|----------------|-------|-------|
| 1   | 10             |       |       |
| 2   | 20             | 10    |       |
| 3   | 30             | 20    | 10    |
| 4   | 40             | 30    | 20    |
| 5   | 50             | 40    | 30    |

The autocorrelation is computed according to the formula:

$$r_k = \frac{\sum_{i=k+1}^n (Y_i - \bar{Y})(Y_{i-k} - \bar{Y})}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \text{ where } k = 0, 1, 2, \dots, n$$

Where  $Y_t$  is the Observed Value at time  $t$ ,  $\bar{Y}$  is the mean of the Observed Values and  $Y_{t-k}$  is the value for Lag- $k$ .

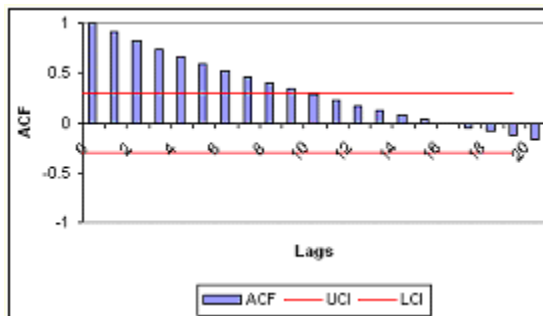
For example, using the values above, the autocorrelation for Lag-1 and Lag - 2 can be calculated as follows.

$$\bar{Y} = (10 + 20 + 30 + 40 + 50) / 5 = 30$$

$$r_1 = ((20 - 30) * (10 - 30) + (30 - 30) * (20 - 30) + (40 - 30) * (30 - 30) + (50 - 30) * (40 - 30)) / ((10 - 30)^2 + (20 - 30)^2 + (30 - 30)^2 + (40 - 30)^2 + (50 - 30)^2) = 0.4$$

$$r_2 = ((30 - 30) * (10 - 30) + (40 - 30) * (20 - 30) + (50 - 30) * (30 - 30)) / (((10 - 30)^2 + (20 - 30)^2 + (30 - 30)^2 + (40 - 30)^2 + (50 - 30)^2) = -0.1$$

The two red horizontal lines on the graph below delineate the Upper confidence level (UCL) and the Lower confidence level (LCL). If the data is random, then the plot should be within the UCL and LCL. If the plot exceeds either of these two levels, as seen in the plot above, then it can be presumed that some correlation exists in the data.



## Partial Autocorrelation Function (PACF)

This technique is used to compute and plot the partial autocorrelations between the original series and the lags. However, PACF eliminates all linear dependence in the time series beyond the specified lag.

## ARIMA

An ARIMA (autoregressive integrated moving-average models) model is a regression-type model that includes autocorrelation. The basic assumption in estimating the ARIMA coefficients is that the data are stationary, that is, the

trend or seasonality cannot affect the variance. This is generally not true. To achieve the stationary data, Analytic Solver Data Science will first apply “differencing”: ordinary, seasonal or both.

After Analytic Solver Data Science fits the model, various results will be available. The quality of the model can be evaluated by comparing the time plot of the actual values with the forecasted values. If both curves are close, then it can be assumed that the model is a good fit. The model should expose any trends and seasonality, if any exist. If the residuals are random then the model can be assumed a good fit. However, if the residuals exhibit a trend, then the model should be refined. Fitting an ARIMA model with parameters (0,1,1) will give the same results as exponential smoothing. Fitting an ARIMA model with parameters (0,2,2) will give the same results as double exponential smoothing.

## Partitioning

To avoid over fitting of the data and to be able to evaluate the predictive performance of the model on new data, we must first partition the data into training and validation sets using Analytic Solver Data Science’s time series partitioning utility. After the data is partitioned, ACF, PACF, and ARIMA can be applied to the dataset.

---

## Examples for Time Series Analysis

The examples below illustrate how Analytic Solver Data Science can be used to explore the Income.xlsx dataset to uncover trends and seasonalities in a dataset. Click **Help – Examples** on the Data Science ribbon, then **Forecasting/Data Science Examples** and open the example dataset, **Income.xlsx**. This dataset contains the average income of tax payers by state.

Typically the following steps are performed in a time series analysis.

1. The data is first partitioned into two sets with 60% of the data assigned to the training set and 40% of the data assigned to validation.
2. Exploratory techniques are applied to both the training and validation sets. If the results are in synch then the model can be fit. If the ACF and PACF plots are the same, then the same model can be used for both sets.
3. The model is fit using the ARIMA method.
4. When we fit a model using the ARIMA method, Analytic Solver displays the ACF and PACF plots for residuals. If these plots are in the band of UCL and LCL then it indicates that the residuals are random and the model is adequate.
2. If the residuals are not within the bands, then some correlation exists, and the model should be improved.

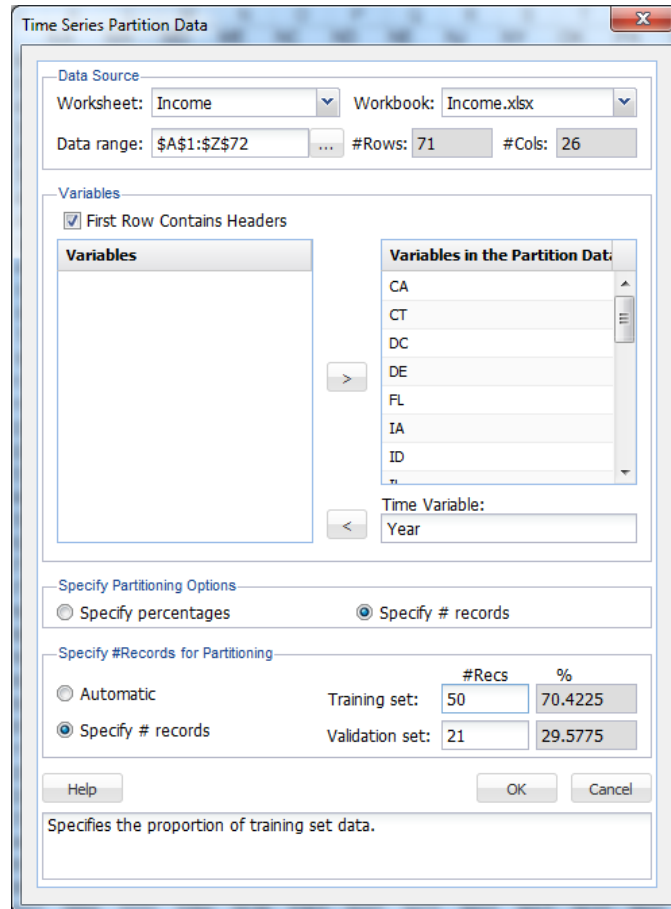
First we must perform a partition on the data. Click **Partition** within the Time Series group on the Data Science ribbon to open the following dialog.

Select **Year** under *Variables* and click > to define the variable as the *Time Variable*. Select the remaining variables under *Variables* and click > to include them in the partitioned data.

Select **Specify #Records** under *Specify Partitioning Options* to specify the number of records assigned to the training and validation sets. Then select

**Specify #Records** under *Specify #Records for Partitioning*. Enter **50** for the number of *Training Set* records and **21** for the number of *Validation Set* records.

If **Specify Percentages** is selected under *Specify Partitioning Options*, Analytic Solver Data Science will assign a percentage of records to each set according to the values entered by the user or automatically entered by Analytic Solver Data Science under *Specify Percentages for Partitioning*.



Click **OK**. *TSPartition* is inserted into the Model tab of the Solver task pane under Reports – Time Series Partition – Run 1.

| Partition Summary |           |
|-------------------|-----------|
| Partition         | # Records |
| Training          | 50        |
| Validation        | 21        |

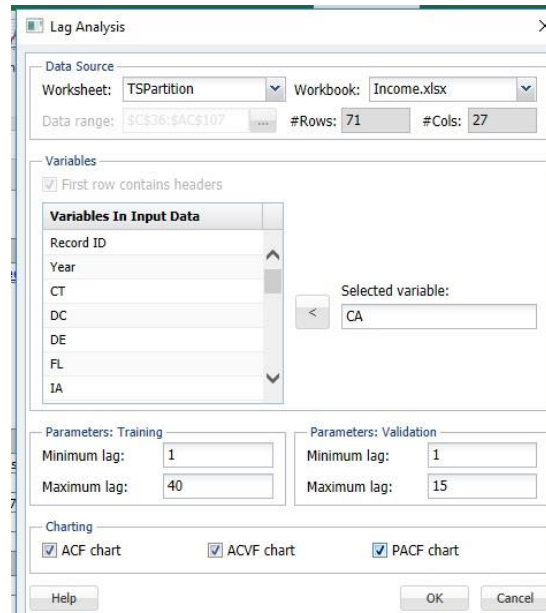
  

| Record ID | Year | CA  | CT   | DC   | DE   | FL  | IA  | ID  | IL  | IN  | KA  | MA  | MD  | ME  | NC  | ND  | NE  | NJ |
|-----------|------|-----|------|------|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|----|
| Record 1  | 1929 | 991 | 1024 | 1269 | 1032 | 518 | 581 | 507 | 948 | 607 | 332 | 906 | 768 | 601 | 332 | 382 | 868 | 9  |
| Record 2  | 1930 | 887 | 921  | 1266 | 857  | 470 | 510 | 553 | 837 | 514 | 467 | 836 | 712 | 576 | 292 | 311 | 833 | 8  |
| Record 3  | 1931 | 749 | 801  | 1198 | 775  | 398 | 400 | 374 | 671 | 438 | 401 | 759 | 638 | 491 | 248 | 187 | 652 | 7  |
| Record 4  | 1932 | 580 | 620  | 1054 | 560  | 319 | 297 | 274 | 486 | 310 | 266 | 613 | 512 | 377 | 187 | 176 | 550 | 5  |
| Record 5  | 1933 | 546 | 583  | 903  | 564  | 288 | 253 | 227 | 437 | 294 | 230 | 559 | 468 | 371 | 208 | 148 | 495 | 5  |

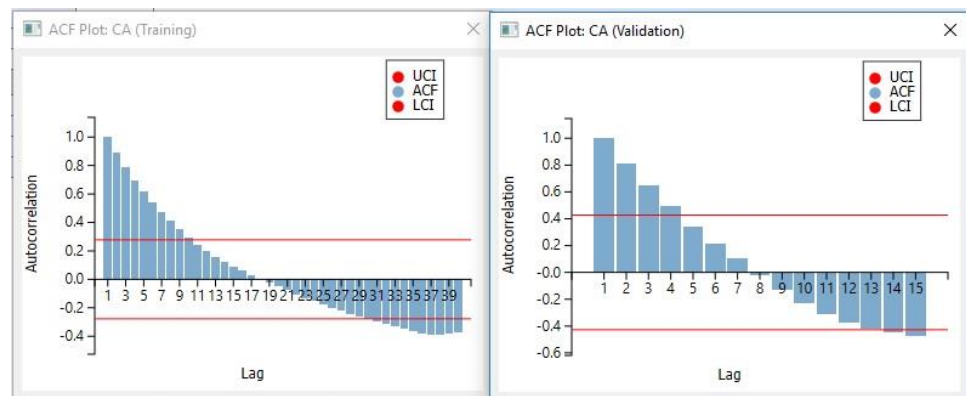
Note in the output above, the partitioning method is sequential (rather than random). The first 50 observations have been assigned to the training set and the remaining 21 observations have been assigned to the validation set.

Open the Lag Analysis dialog by clicking **ARIMA – Lag Analysis**. Select **CA** under *Variables In input data*, then click > to move the variable to *Selected variable*. Enter **1** for Minimum Lag and **40** for *Maximum Lag* under *Parameters: Training* and 1 for Minimum Lag and 15 for *Maximum Lag* under *Parameters: Validation*.

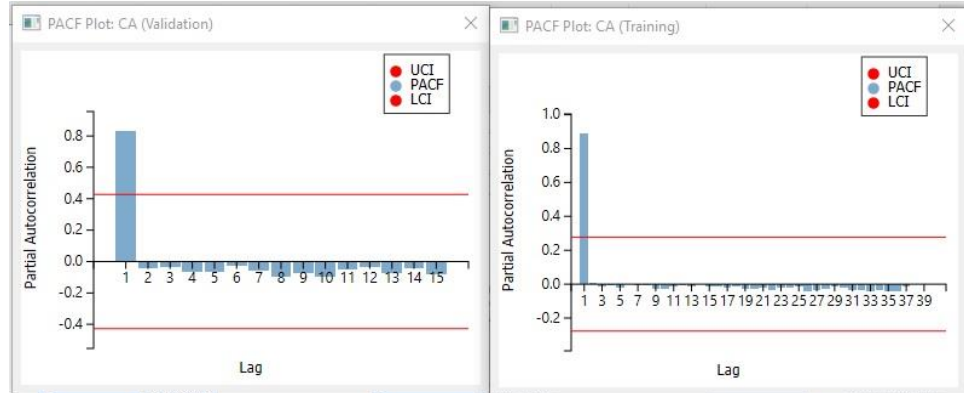
Under Charting, select ACF chart, ACVF chart, and PACF chart to include each chart in the output.



Click **OK**. *TS\_Lags* is inserted into the task pane under Reports – Autocorrelations – Run 1.

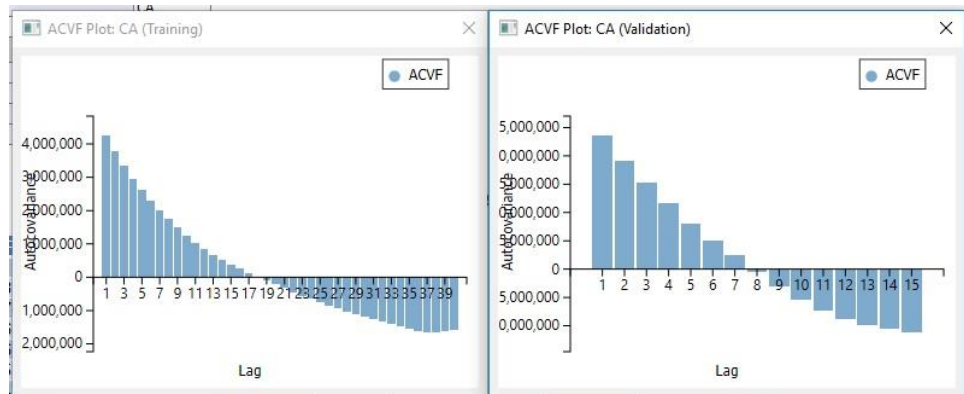


First, let's take a look at the ACF charts. Note on each chart, the autocorrelation decreases as the number of lags increase. This suggests that a definite pattern does exist in each partition. However, since the pattern does not repeat, it can be assumed that no seasonality is included in the data. In addition, both charts appear to exhibit a similar pattern.



The PACF functions show a definite pattern which means there is a trend in the data. However, since the pattern does not repeat, we can conclude that the data does not show any seasonality.

A plot of the autocovariance values has been added to the output.



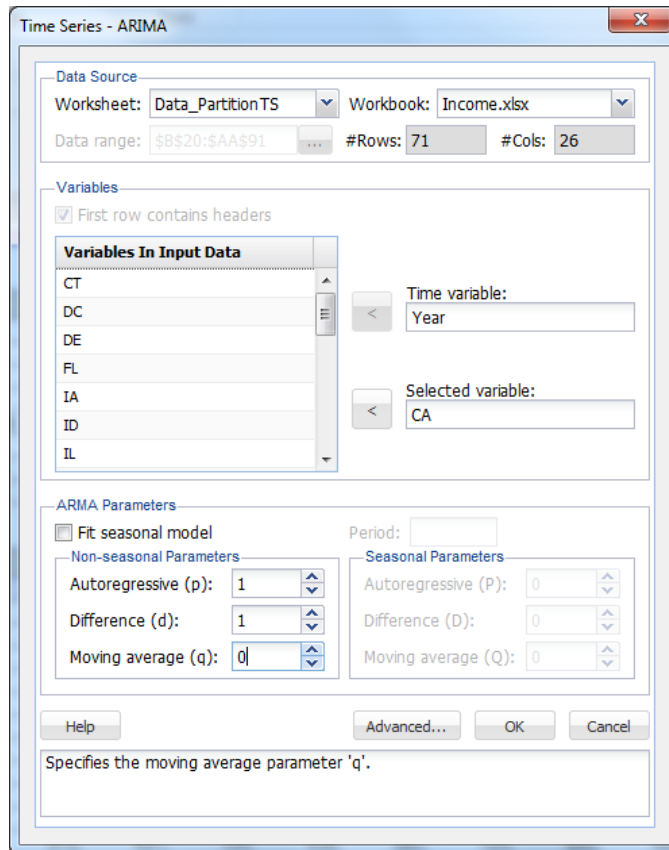
All three charts suggest that a definite pattern exists in the data, but no seasonality. In addition, both datasets exhibit the same behavior in both the training and validation sets which suggests that the same model could be appropriate for each. Now we are ready to fit the model.

The ARIMA model accepts three parameters:  $p$  – the number of autoregressive terms,  $d$  – the number of non-seasonal differences, and  $q$  – the number of lagged errors (moving averages).

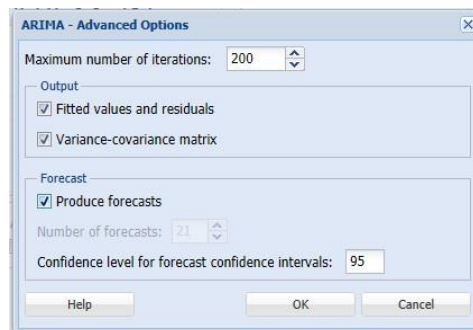
Recall that the ACF plot showed no seasonality in the data which means that autocorrelation is almost static, decreasing with the number of lags increasing. This suggests setting  $q = 0$  since there appears to be no lagged errors. The PACF plot displayed a large value for the first lag but minimal plots for successive lags. This suggest setting  $p = 1$ . With most datasets, setting  $d = 1$  is sufficient or can at least be a starting point.

Click back to the TSPartition tab and then click **ARIMA – ARIMA Model** to bring up the *Time Series – ARIMA* dialog.

Select **CA** under *Variables In input data* then click > to move the variable to the *Selected Variable* field. Under *Nonseasonal Parameters* set **Autoregressive (p)** to **1**, **Difference (d)** to **1** and **Moving Average (q)** to **0**.



Click **Advanced** to open the *ARIMA – Advanced Options* dialog. Select **Fitted Values and residuals**, **Produce forecasts**, and **Report Forecast Confidence Intervals**. The default *Confidence Level* setting of 95 is automatically entered. The option *Variance-covariance matrix* is selected by default.



Click **OK** on the *ARIMA-Advanced Options* dialog and again on the *Time Series – ARIMA* dialog. Analytic Solver Data Science calculates and displays various parameters and charts in four output sheets, *Arima\_Output*, *Arima\_Fitted*, *Arima\_Forecast* and *Arima\_Stored*. Click the *Arima\_Output* tab to view the Output Navigator.

|   | A | B | C                       | D                                            | E                                                  | F                                          | G                        | H                                        | I                                                    | J                                          | K                      | L                           | M                           |
|---|---|---|-------------------------|----------------------------------------------|----------------------------------------------------|--------------------------------------------|--------------------------|------------------------------------------|------------------------------------------------------|--------------------------------------------|------------------------|-----------------------------|-----------------------------|
| 2 |   |   |                         |                                              |                                                    |                                            |                          |                                          |                                                      |                                            |                        |                             |                             |
| 3 |   |   | <b>Output Navigator</b> |                                              |                                                    |                                            |                          |                                          |                                                      |                                            |                        |                             |                             |
| 4 |   |   | <a href="#">Fitted</a>  | <a href="#">Ljung-Box Test for Residuals</a> | <a href="#">Lag Analysis: Residuals - Training</a> | <a href="#">Variance-Covariance Matrix</a> | <a href="#">Forecast</a> | <a href="#">Error Measures: Training</a> | <a href="#">Lag Analysis: Residuals - Validation</a> | <a href="#">Error Measures: Validation</a> | <a href="#">Inputs</a> | <a href="#">PMMML Model</a> | <a href="#">ARIMA Model</a> |
| 5 |   |   |                         |                                              |                                                    |                                            |                          |                                          |                                                      |                                            |                        |                             |                             |

Click the *ARIMA Model* link on the Output Navigator to move to display the ARIMA Model and Ljung-Box Test Results on Residuals.



| Record ID | Coeff     | Std-Dev     | p-value   |
|-----------|-----------|-------------|-----------|
| Const     | -16.20163 | 3.053324388 | 1.119E-07 |
| AR 1      | 1.0920227 | 0.054394554 | 1.198E-89 |

|             |           |
|-------------|-----------|
| Mean        | 176.06122 |
| -2LogL      | 592.00842 |
| Res. StdDev | 103.97445 |
| #Iterations | 7         |

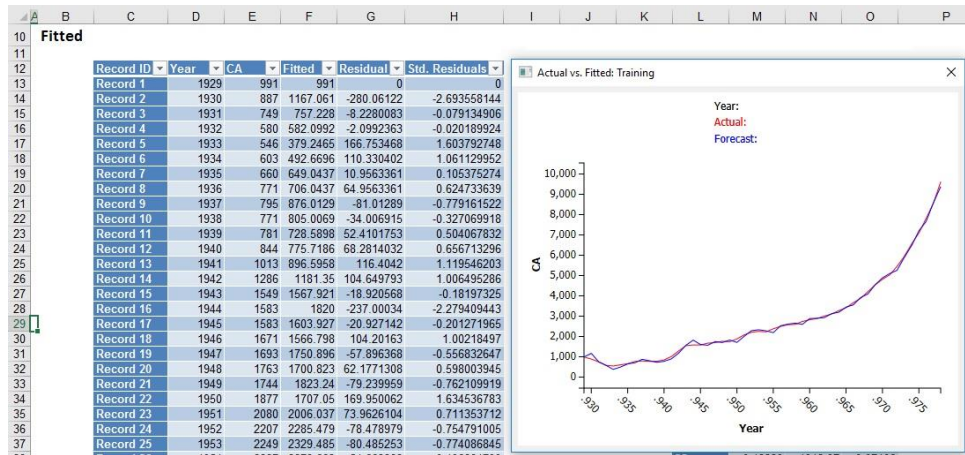
  

| Record ID | p-value   | chi-square  | df |
|-----------|-----------|-------------|----|
| Lag 12    | 0.9999975 | 14.34135727 | 11 |
| Lag 24    | 1         | 30.77988978 | 23 |
| Lag 36    | 1         | 48.16816013 | 35 |
| Lag 48    | 1         | 85.16075414 | 47 |

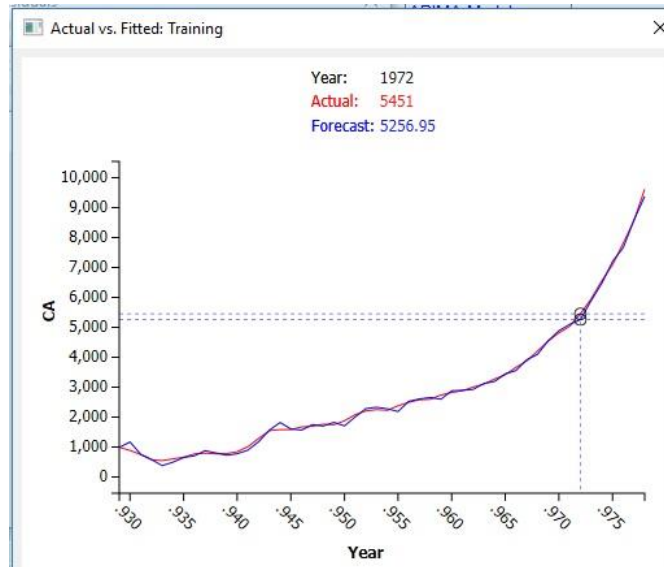
Analytic Solver has calculated the constant term and the AR1 term for our model, as seen above. These are the constant and f1 terms of our forecasting equation. See the following output of the Chi - square test.

The very small p-values for the constant term (1.119E-7) and AR1 term (1.19e-89) suggest that the model is a good fit to our data.

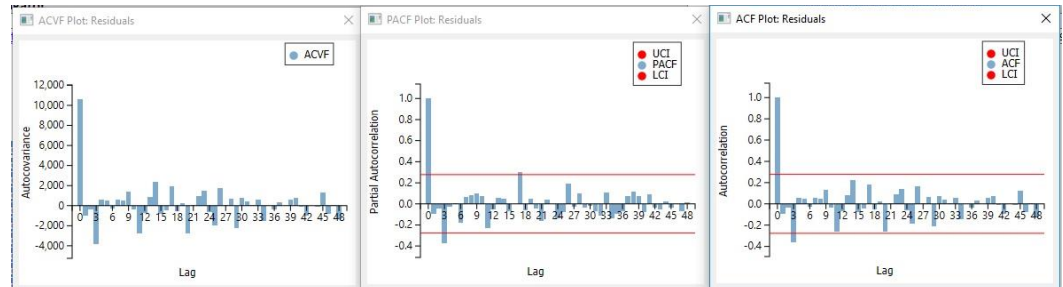
Click the Lag Analysis: Residuals – Training link. This table plots the actual and fitted values and the resulting residuals for the training partition. As shown in the graph below, the Actual and Forecasted values match up fairly well. The usefulness of the model in forecasting will depend upon how close the actual and forecasted values are in the Forecast, which we will inspect later.



Use your mouse to select a point on the graph to compare the Actual value to the Forecasted value.

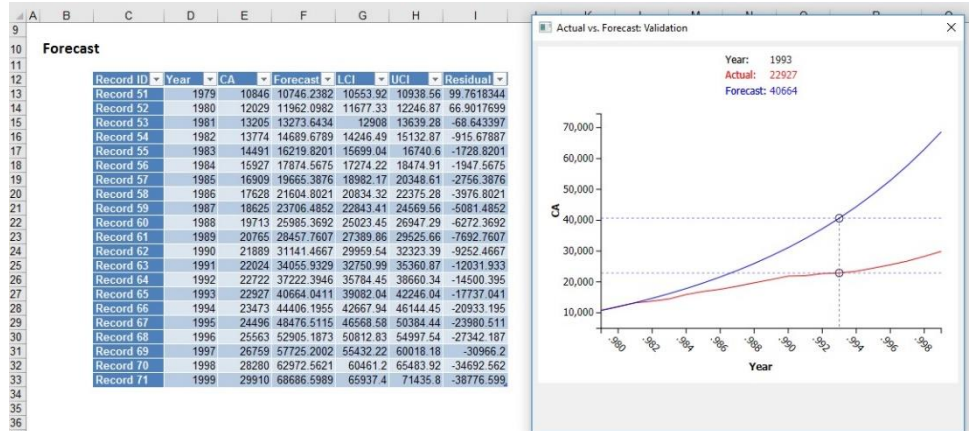


Take a look at the ACF and PACF plots for Errors found at the bottom of *ARIMA\_Output*. One additional chart was added starting in V2017 - the ACVF Plot for the Residuals.



With the exception of Lag 1, the majority of the lags in the PACF and ACF charts are either clearly within the UCL and LCL bands or just outside of these bands. This suggests that the residuals are random and are not correlated.

Click the Forecast link on the Output Navigator to display the Forecast Data table and charts.



The table shows the actual and forecasted values along with LCI (Lower Confidence Interval), UCI (Upper Confidence Interval) and Residual values. The "Lower" and "Upper" values represent the lower and upper bounds of the confidence interval. There is a 95% chance that the forecasted value will fall into this range. The graph to the right plots the Actual values for CA against the Forecasted values. Again, click any point on either curve to compare the Actual against the Forecasted values.

# Automated Risk Analysis of Machine Learning Models

---

## Introduction

This example illustrates the use of Analytic Solver Data Science to assess the uncertainty and risk of a machine learning (ML) model. The model is “trained” and “validated” on a total of 9,578 known cases of loan applicants, including “outcome” data on whether they defaulted or were fully repaid. Before the model is put into “production use”, Analytic Solver Data Science is used to perform, in a fully automated manner, a risk analysis of the model’s performance on 5,747 *new* cases which are individually different, but whose overall behavior is consistent with the known data. The risk analysis also includes a quantitative assessment of the range of possible loan losses, if loan decisions are made on the basis of this trained ML model.

---

## Risk Analysis of ML Model in Predicting Loan Defaults

The dataset used in this example originated from LendingClub.com and is publicly available on the Kaggle machine learning website at <https://www.kaggle.com/datasets/itsuru/loan-data>. The “features” or independent variables in this dataset are summarized below:

1. `credit.policy`: 1 if the customer meets the credit underwriting criteria of LendingClub.com, and 0 otherwise.
2. `purpose`: The purpose of the loan (takes values "creditcard", "debtconsolidation", "educational", "majorpurchase", "smallbusiness", and "all\_other").
3. `int.rate`: The interest rate of the loan, as a proportion (a rate of 11% would be stored as 0.11). Borrowers judged by LendingClub.com to be more risky are assigned higher interest rates.
4. `installment`: The monthly installments owed by the borrower if the loan is funded.
5. `log.annual.inc`: The natural log of the self-reported annual income of the borrower.
6. `dti`: The debt-to-income ratio of the borrower (amount of debt divided by annual income).
7. `fico`: The FICO credit score of the borrower.
8. `days.with.cr.line`: The number of days the borrower has had a credit line.
9. `revol.bal`: The borrower's revolving balance (amount unpaid at the end of the credit card billing cycle).
10. `revol.util`: The borrower's revolving line utilization rate (the amount of the credit line used relative to total credit available).
11. `inq.last.6mths`: The borrower's number of inquiries by creditors in the last 6 months.
12. `delinq.2yrs`: The number of times the borrower had been 30+ days past due on a payment in the past 2 years.

- 13. pub.rec: The borrower's number of derogatory public records (bankruptcy filings, tax liens, or judgments).
- 14. not.fully.paid: 1 if the borrower defaulted and did not [fully] pay off the loan.

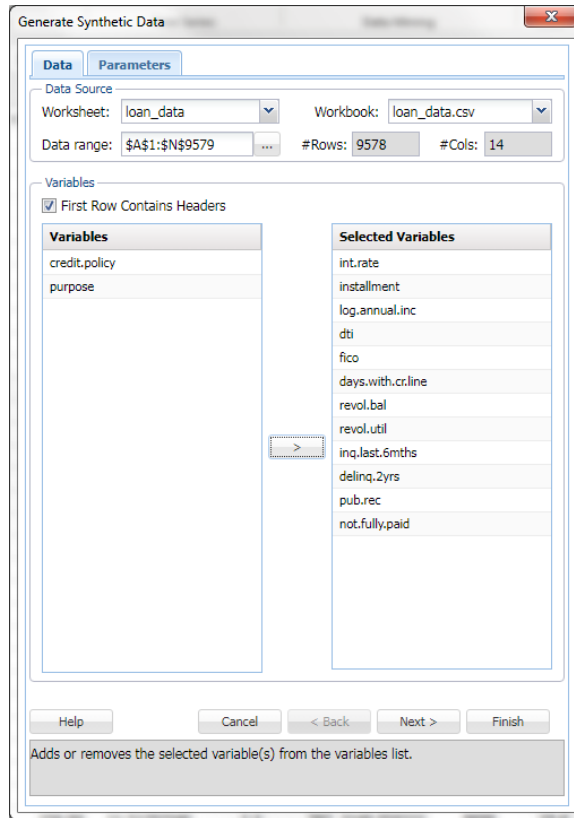
In this example, the risk analysis reveals that the first machine learning model, using the Classification Tree ML method, has a high risk of poor performance on future cases. As a result a second ML model is fit using the Logistic Regression ML model. In this case the risk analysis shows lower risk of poor results, but still yields a *quantification* of that *risk*, not available before in Analytic Solver Data Science.

The first figure shows the loan dataset loaded into Microsoft Excel, for use with Analytic Solver Data Science for Excel.

|    | A             | B                  | C        | D           | E              | F     | G    | H            | I         | J          | K           | L          | M       | N              | O |
|----|---------------|--------------------|----------|-------------|----------------|-------|------|--------------|-----------|------------|-------------|------------|---------|----------------|---|
|    | credit.policy | purpose            | int.rate | installment | log.annual.inc | dti   | fico | days.with.cr | revol.bal | revol.util | inq.last.6r | delinq.2yr | pub.rec | not.fully.paid |   |
| 1  |               |                    |          |             |                |       |      |              |           |            |             |            |         |                |   |
| 2  | 1             | debt_consolidation | 0.1189   | 829.1       | 11.35040554    | 19.48 | 727  | 5639.958333  | 28854     | 52.1       | 0           | 0          | 0       | 0              |   |
| 3  | 1             | credit_card        | 0.1071   | 228.22      | 11.08214255    | 14.29 | 707  | 2760         | 33623     | 76.7       | 0           | 0          | 0       | 0              |   |
| 4  | 1             | debt_consolidation | 0.1357   | 366.86      | 10.37349118    | 11.63 | 682  | 4710         | 3511      | 25.6       | 1           | 0          | 0       | 0              |   |
| 5  | 1             | debt_consolidation | 0.1008   | 162.34      | 11.35040654    | 8.1   | 712  | 2699.958333  | 33667     | 73.2       | 1           | 0          | 0       | 0              |   |
| 6  | 1             | credit_card        | 0.1426   | 102.92      | 11.29973224    | 14.97 | 667  | 4066         | 4740      | 39.5       | 0           | 1          | 0       | 0              |   |
| 7  | 1             | credit_card        | 0.0788   | 125.13      | 11.90436755    | 16.98 | 727  | 6120.041667  | 50807     | 51         | 0           | 0          | 0       | 0              |   |
| 8  | 1             | debt_consolidation | 0.1496   | 194.02      | 10.71441777    | 4     | 667  | 3180.041667  | 3839      | 76.8       | 0           | 0          | 1       | 1              |   |
| 9  | 1             | all_other          | 0.1114   | 131.22      | 11.60209984    | 11.16 | 727  | 5559.958333  | 24220     | 68.6       | 0           | 0          | 0       | 1              |   |
| 10 | 1             | home_improvement   | 0.1134   | 87.19       | 11.40756495    | 17.25 | 682  | 3989         | 69909     | 51.1       | 1           | 0          | 0       | 0              |   |
| 11 | 1             | debt_consolidation | 0.1221   | 84.12       | 10.20359214    | 10    | 707  | 2730.041667  | 5630      | 23         | 1           | 0          | 0       | 0              |   |
| 12 | 1             | debt_consolidation | 0.1347   | 360.43      | 10.4341158     | 22.09 | 677  | 6713.041667  | 13846     | 71         | 2           | 0          | 1       | 0              |   |
| 13 | 1             | debt_consolidation | 0.1124   | 253.58      | 11.83500896    | 9.16  | 662  | 4298         | 5122      | 18.2       | 2           | 1          | 0       | 0              |   |
| 14 | 1             | debt_consolidation | 0.0859   | 316.11      | 10.93310697    | 15.49 | 767  | 6519.958333  | 6068      | 16.7       | 0           | 0          | 0       | 0              |   |
| 15 | 1             | small_business     | 0.0714   | 92.82       | 11.51292546    | 6.5   | 747  | 4384         | 3021      | 4.8        | 0           | 1          | 0       | 0              |   |
| 16 | 1             | debt_consolidation | 0.0863   | 209.54      | 9.48197269     | 9.73  | 727  | 1559.958333  | 6282      | 44.6       | 0           | 0          | 0       | 0              |   |
| 17 | 1             | major_purchase     | 0.1103   | 327.53      | 10.75891524    | 13.04 | 702  | 8159.958333  | 5394      | 53.4       | 1           | 0          | 0       | 0              |   |
| 18 | 1             | all_other          | 0.1317   | 77.69       | 10.52277288    | 2.26  | 672  | 3895.958333  | 2211      | 88.4       | 0           | 0          | 0       | 0              |   |
| 19 | 1             | credit_card        | 0.0894   | 476.58      | 11.60823564    | 7.07  | 797  | 6510.958333  | 7586      | 52.7       | 1           | 0          | 0       | 0              |   |
| 20 | 1             | debt_consolidation | 0.1039   | 584.12      | 10.49127422    | 3.8   | 712  | 2760         | 8311      | 59.8       | 0           | 0          | 0       | 0              |   |
| 21 | 1             | major_purchase     | 0.1513   | 173.65      | 11.00209984    | 2.74  | 667  | 1126.958333  | 591       | 84.4       | 3           | 0          | 0       | 0              |   |
| 22 | 1             | all_other          | 0.08     | 188.02      | 11.22524339    | 16.08 | 772  | 4888.958333  | 29797     | 23.2       | 1           | 0          | 0       | 0              |   |
| 23 | 1             | all_other          | 0.0863   | 474.42      | 10.81977828    | 2.59  | 797  | 11951        | 5656      | 27.6       | 0           | 0          | 0       | 0              |   |
| 24 | 1             | credit_card        | 0.1355   | 339.6       | 11.51292546    | 7.94  | 662  | 1939.958333  | 21162     | 57.7       | 0           | 0          | 0       | 0              |   |
| 25 | 1             | credit_card        | 0.0788   | 484.85      | 11.73606902    | 7.05  | 782  | 5640.041667  | 16931     | 34.6       | 1           | 0          | 0       | 0              |   |
| 26 | 1             | debt_consolidation | 0.1229   | 320.19      | 11.26446411    | 8.8   | 672  | 3760.958333  | 4822      | 58.1       | 0           | 0          | 1       | 0              |   |
| 27 | 1             | all_other          | 0.0901   | 159.03      | 12.4292162     | 10    | 712  | 1553.958333  | 14354     | 36.6       | 0           | 2          | 0       | 0              |   |
| 28 | 1             | all_other          | 0.0743   | 155.38      | 11.08214255    | 0.28  | 802  | 4649.958333  | 1576      | 5.7        | 1           | 0          | 0       | 0              |   |
| 29 | 1             | debt_consolidation | 0.1175   | 295.43      | 9.99879732     | 14.29 | 662  | 1318.958333  | 4175      | 51.5       | 0           | 1          | 0       | 0              |   |
| 30 | 1             | all_other          | 0.0743   | 155.38      | 12.20607265    | 0.28  | 772  | 4516.958333  | 3164      | 13.7       | 0           | 0          | 0       | 0              |   |
| 31 | 1             | all_other          | 0.0743   | 155.38      | 12.20607265    | 3.72  | 812  | 6778.958333  | 85607     | 0.7        | 0           | 0          | 0       | 0              |   |
| 32 | 1             | debt_consolidation | 0.0807   | 156.84      | 11.51292546    | 2.3   | 742  | 3148.958333  | 9698      | 19.4       | 0           | 0          | 0       | 0              |   |
| 33 | 1             | credit_card        | 0.1028   | 275.38      | 9.798127037    | 6.4   | 692  | 7469.958333  | 8847      | 26.9       | 1           | 1          | 0       | 0              |   |
| 34 | 1             | home_improvement   | 0.0743   | 155.38      | 11.91839057    | 0     | 777  | 7128.958333  | 6053      | 19.5       | 0           | 0          | 0       | 0              |   |
| 35 | 1             | home_improvement   | 0.0807   | 78.42       | 11.60823564    | 11.33 | 762  | 6059.958333  | 7274      | 13.1       | 0           | 0          | 0       | 0              |   |
| 36 | 1             | credit_card        | 0.087    | 158.3       | 11.22524339    | 15.55 | 757  | 4770         | 66033     | 23         | 0           | 0          | 0       | 0              |   |

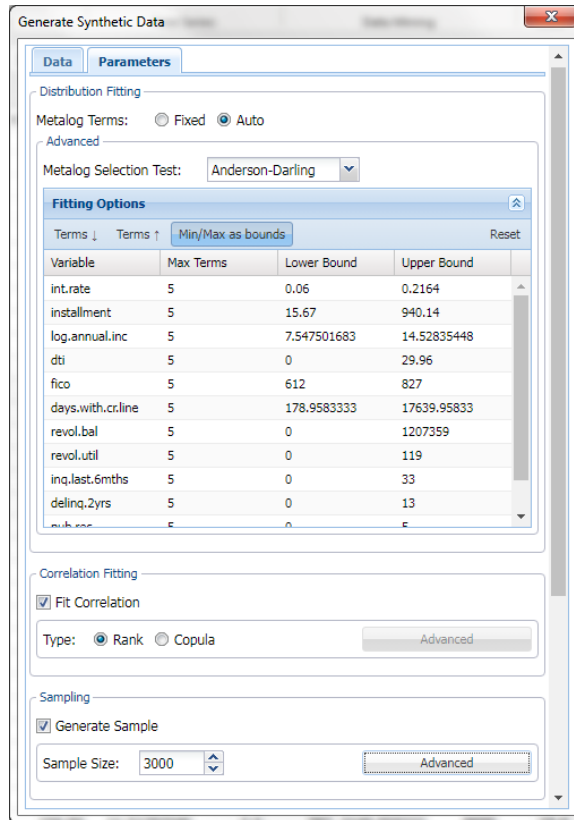
## Preliminary Illustration to Show Synthetic Data Compared to Known Data

We select the 12 numeric/continuous “features” in this dataset, omitting the “credit.policy” and “purpose” (categorical and text) features, and click the “Next” button at the bottom of the dialog:



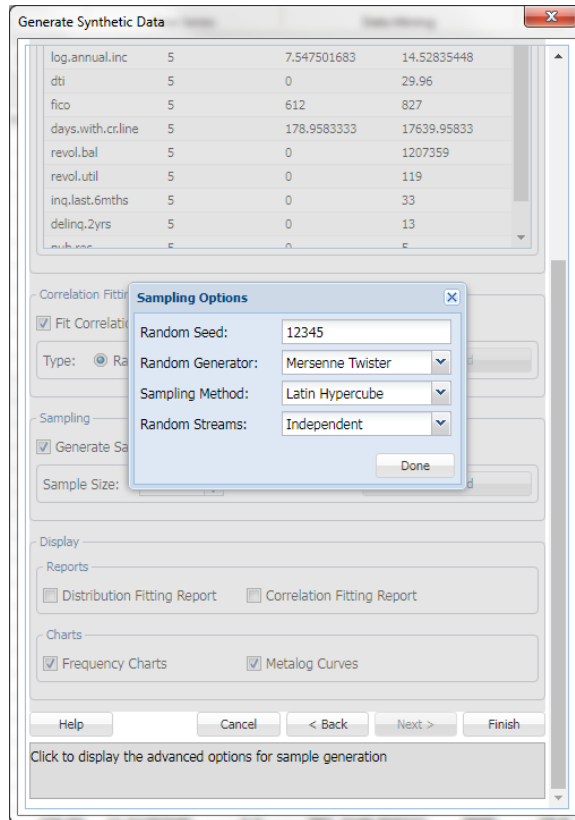
The next two images, from the same step “Parameters” dialog, illustrate settings that the user can *optionally* modify to further control the simulation and synthetic data generation process; however, in most applications the user need not change anything here, and can simply click the “Next” or “Finish” buttons.

For more information on all Classification Tree options, see the Classification Tree chapter within the Data Science Reference Guide.



Above, we’ve expanded the “Fitting Options”, which control the automated fitting of Metalog probability distributions to each feature. For example, by default, bounds in the dataset are used as bounds for the family of Metalog distributions, but users can easily and quickly adjust or remove the lower and/or upper bounds. “Correlation Fitting” will by default construct a rank-order correlation matrix that includes all of the features, but choosing “Copula” activates the “Advanced” button, where types of copulas (Clayton, Frank, Gumbel, etc.) can be chosen.

For more information on these options, see the Generate Data chapter found within the Data Science Reference Guide.



Above, we’re illustrating the lower portion of the same “Parameters” dialog, with further options including the “Sampling Options” shown. Again, in typical applications the user can simply click the “Next” or “Finish” buttons, without choosing anything here.

The next image shows frequency histograms of the *known* cases versus the *simulated* or synthetically generated cases: It is the *differences* (rather than the similarities) in this data that give rise to the uncertainty and risk in the performance of the machine learning model.





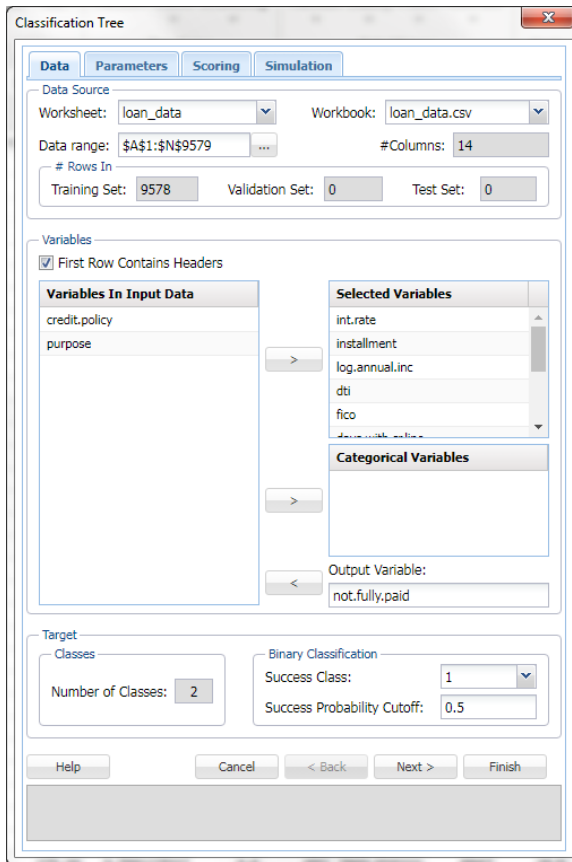
In the above frequency histograms, the original data (known cases) appear in the top row, and the simulated or synthetically generated cases appear in the bottom row. Note that the distribution of “dti” (debt to income ratio) is *different* in the simulated cases, with more frequent cases having a higher debt to income ratio.

## Illustration: Automated Risk Analysis of Just-Trained Machine Learning Model

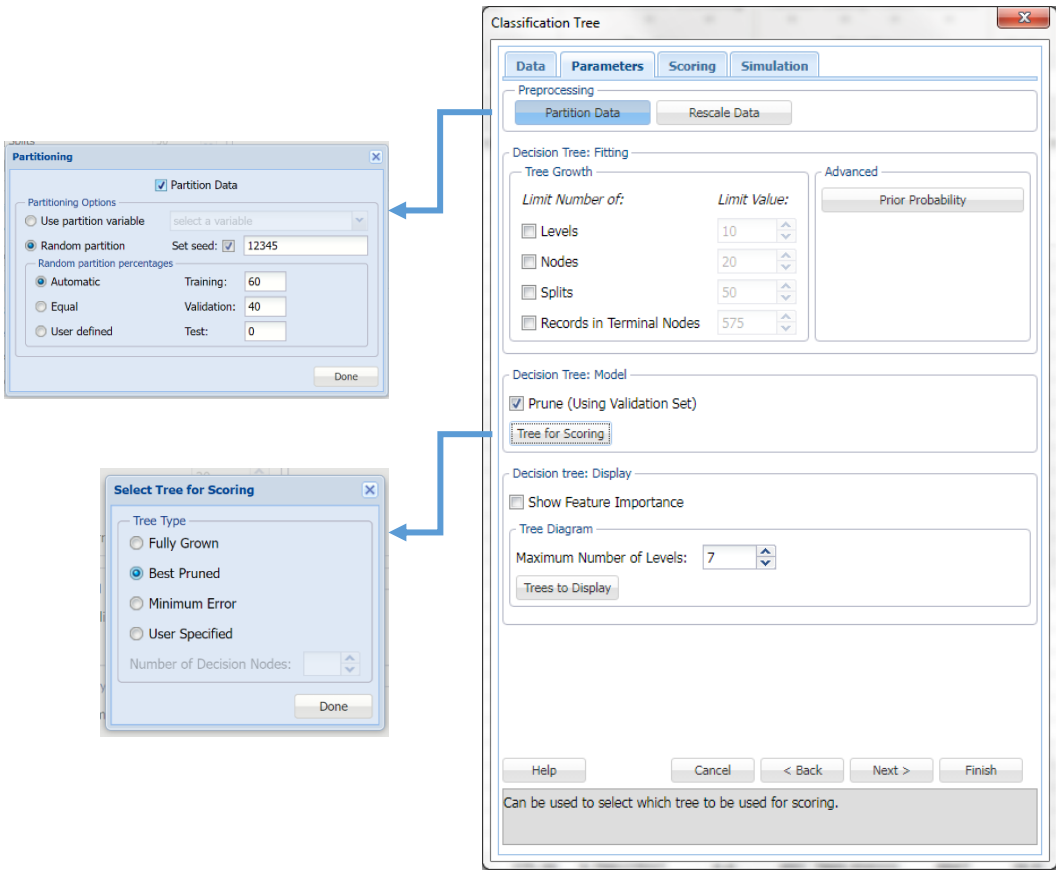
With the above as background, the fully automated risk analysis process is illustrated. First, a “Classification Tree”, a typical ML model based on the CART methodology is built:

| credit.policy | purpose              | intr.rate | installment | log.annual.inc | dti   | fico | dti         |
|---------------|----------------------|-----------|-------------|----------------|-------|------|-------------|
| 1             | 1 debt_consolidation | 0.1189    | 428.11      | 11.35040654    | 19.48 | 737  | 5           |
| 2             | 1 credit_card        | 0.1071    | 228.22      | 11.08214255    | 14.29 | 707  | 5           |
| 3             | 1 debt_consolidation | 0.1357    | 366.86      | 10.37349118    | 11.63 | 682  | 2           |
| 4             | 1 debt_consolidation | 0.1008    | 162.34      | 11.35040654    | 8.1   | 712  | 2           |
| 5             | 1 credit_card        | 0.1426    | 102.92      | 11.29973224    | 14.97 | 667  | 6           |
| 6             | 1 credit_card        | 0.0788    | 125.13      | 11.90496755    | 16.98 | 727  | 6           |
| 7             | 1 debt_consolidation | 0.1496    | 194.02      | 10.71441777    | 4     | 667  | 3           |
| 8             | 1 all_other          | 0.1114    | 131.22      | 11.00209984    | 11.08 | 722  | 2           |
| 9             | 1 home_improvement   | 0.1134    | 371.59      | 11.40756495    | 17.25 | 682  | 2           |
| 10            | 1 debt_consolidation | 0.1221    | 84.12       | 10.20359214    | 10    | 707  | 2           |
| 11            | 1 debt_consolidation | 0.1347    | 360.43      | 10.4341158     | 22.09 | 677  | 6           |
| 12            | 1 debt_consolidation | 0.1324    | 253.58      | 11.83500896    | 9.16  | 662  | 6           |
| 13            | 1 debt_consolidation | 0.0859    | 316.11      | 10.93310697    | 15.49 | 767  | 6           |
| 14            | 1 small_business     | 0.0714    | 92.82       | 11.51292546    | 6.5   | 747  | 6           |
| 15            | 1 debt_consolidation | 0.0863    | 209.54      | 9.487972109    | 9.73  | 727  | 1559.958333 |
| 16            | 1 major_purchase     | 0.1103    | 327.53      | 10.73891524    | 13.04 | 702  | 8159.958333 |
| 17            | 1 all_other          | 0.1317    | 77.69       | 10.3227288     | 2.26  | 672  | 3895.958333 |
| 18            | 1 credit_card        | 0.0894    | 476.58      | 11.60823564    | 7.07  | 797  | 6510.958333 |
| 19            | 1 debt_consolidation | 0.1039    | 584.12      | 10.49127422    | 3.8   | 712  | 2760        |
| 20            | 1 major_purchase     | 0.1513    | 173.65      | 11.00209984    | 2.74  | 667  | 1126.958333 |
| 21            | 1 all_other          | 0.08      | 188.02      | 11.22524339    | 16.08 | 772  | 4888.958333 |
| 22            | 1 all_other          | 0.0863    | 474.42      | 10.81977828    | 2.59  | 797  | 11951       |
| 23            | 1 credit_card        | 0.1355    | 339.6       | 11.51292546    | 7.94  | 662  | 1939.958333 |
| 24            | 1 credit_card        | 0.0788    | 484.85      | 11.73609992    | 7.05  | 782  | 5640.041667 |
| 25            | 1 debt_consolidation | 0.1229    | 320.19      | 11.26464411    | 8.8   | 672  | 3760.958333 |
| 26            | 1 all_other          | 0.0901    | 159.03      | 12.4292162     | 10    | 712  | 1553.958333 |
| 27            | 1 all_other          | 0.0743    | 155.38      | 11.08214255    | 0.28  | 802  | 4649.958333 |
| 28            | 1 debt_consolidation | 0.1375    | 255.43      | 9.99879732     | 14.29 | 662  | 1318.958333 |
| 29            | 1 all_other          | 0.0743    | 155.38      | 12.20607265    | 0.28  | 772  | 4516.958333 |
| 30            | 1 all_other          | 0.0743    | 155.38      | 12.20607265    | 3.72  | 812  | 6778.958333 |
| 31            | 1 debt_consolidation | 0.0807    | 156.84      | 11.51292546    | 2.3   | 742  | 3148.958333 |
| 32            | 1 credit_card        | 0.1028    | 275.38      | 9.798127037    | 6.4   | 692  | 7469.958333 |
| 33            | 1 home_improvement   | 0.0743    | 155.38      | 11.91839057    | 0     | 777  | 7128.958333 |
| 34            | 1 home_improvement   | 0.0807    | 78.42       | 11.60823564    | 11.33 | 762  | 6053.958333 |

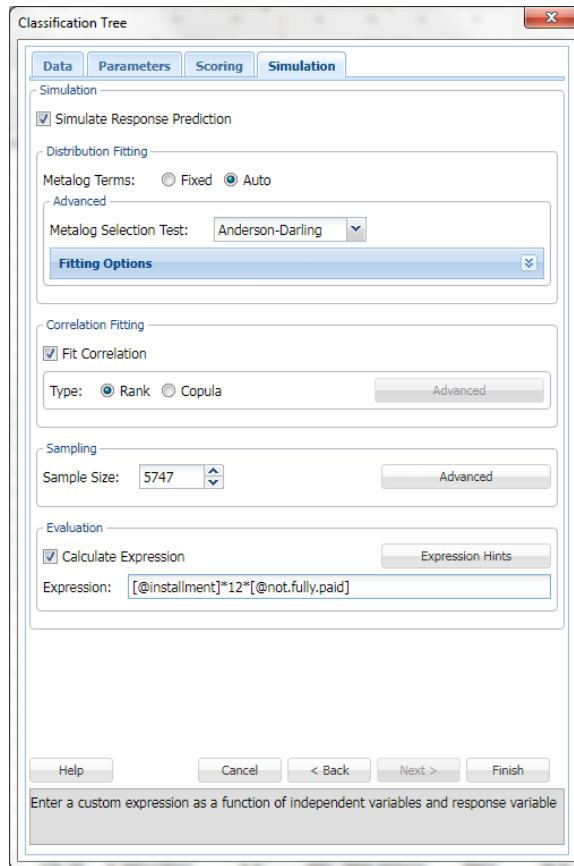
In the first step, the features to be used are selected.



Click the “Next” button to select options for a Classification Tree.



The 9,578 row dataset is first partitioned into a “Training Set” with 60% of the data (5,747 rows), and a “Validation Set” with 40% of the data by clicking Partition Data in the Preprocessing section of the Parameters tab. Select “Prune (Using the Validation Set)”, and then click the Tree for Scoring button to select “Best Pruned” to use the Pruned Tree for Scoring. Then click the “Next” button twice to advance to the Simulation tab. No options will be changed on the Scoring tab.)



Options for the *entire risk analysis* of the *just-trained* and validated machine learning model *“on the fly”* are set on the Simulation tab (shown above). This example asks for 5,747 simulated cases (the same number as in the Training Set, for convenience) and enters a “financial payoff” calculation – via commonly-used Microsoft Excel Table Expression syntax, which calculates a loss of 12 times the monthly installment if the borrower defaults on the loan. Results will appear within moments after clicking the “Finish” button.

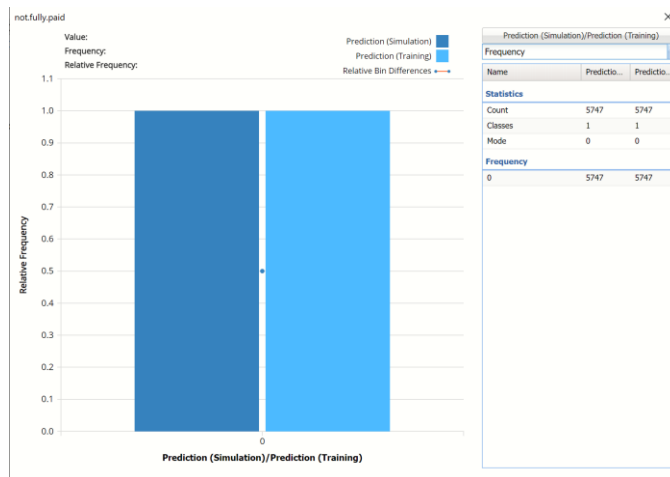
Analytic Solver Data Science now performs the “training” and “validation” steps to create a Classification Tree ML model with parameters fitted to best predict or classify the training set known cases, then validates and refines the model based on the validation set known cases. The trained model is rendered in PMML (Predictive Modeling Markup Language, an industry standard across different software tools).

```

10 PMML Model
11
12 <?xml version="1.0" encoding="utf-8"?>
13 <PMML version="4.2" xmlns:schemaLocation="http://www.dmg.org/PMML-4_2 http://www.dmg.org/v4-2/pmml-4-2.xsd" xmlns:xsi="http://www.w3.org/2001/XMLSchema
14 <Header copyright="Copyright (c) 2022 Frontline Systems Inc." description="TreeModel">
15 <Application name="XLMiner5DX" version="22.0.4.0"/>
16 <Timestamp>2022-8-9 22:14:42</Timestamp>
17 <ModelName/>
18 </Header>
19 <DataDictionary numberOfFields="12">
20 <DataField optype="continuous" dataType="double" name="int.rate"/>
21 <DataField optype="continuous" dataType="double" name="installment"/>
22 <DataField optype="continuous" dataType="double" name="log.annual.inc"/>
23 <DataField optype="continuous" dataType="double" name="dti"/>
24 <DataField optype="continuous" dataType="double" name="fico"/>
25 <DataField optype="continuous" dataType="double" name="days.with.cr.line"/>
26 <DataField optype="continuous" dataType="double" name="revol.bal"/>
27 <DataField optype="continuous" dataType="double" name="revol.util"/>
28 <DataField optype="continuous" dataType="double" name="inq.last.6mths"/>
29 <DataField optype="continuous" dataType="double" name="delinq.2yrs"/>
30 <DataField optype="continuous" dataType="double" name="pub.rec"/>
31 <DataField optype="categorical" dataType="string" name="not.fully.paid">
32 <Value value="0"/>
33 <Value value="1"/>
34 </DataField>
35 </DataDictionary>
36 <TreeModel modelName="TreeModel" functionName="classification" algorithmName="DecisionTree" splitCharacteristic="binarySplit" missingValueStrategy="defau
37 <LocalTransformations/>
38 <MiningSchema>
39 <MiningField name="not.fully.paid" usageType="predicted"/>
40 </MiningSchema>
41 <Output>
42 <OutputField optype="categorical" dataType="string" name="Predicted_not.fully.paid" feature="predictedValue"/>
43 </Output>
44 <Targets>
45 <Target field="not.fully.paid" optype="categorical">

```

The surprise – and financial consequences – appear on the “CT\_Simulation” tab next to the “CT\_Stored” (stored PMML model) tab in the Excel workbook. The simulated data, ML model predictions on that data, and financial consequences are calculated and stored on this “CT\_Simulation” worksheet; clicking that tab displays the chart shown below (in one of its variants).



The light blue bar displays the frequencies of “no default” (0) and “default” (1) on the 5,747 known cases in the Training Set. The dark blue bar displays the same frequencies for the simulated data: This shows that the Classification Tree model is *unlikely* to predict (or classify) *any* future cases as likely to “default”. There is a risk that the model’s future predictive power is not very good at all.

[This is an illustration of an “extreme” case of poor performance, due in part to selecting and training only one type of ML model, which is then “caught” by performing the risk analysis; a more thorough preparation, data exploration, selection and training of alternative models could have yielded better results – but even then, as shown below, the user would lack any quantitative information about the *uncertainty* and *risk* in the finally-chosen ML model.]

In this example, re-examination of the Classification Tree's performance leads the user to try a different statistical / machine learning approach, Logistic Regression:

| credit_policy | purpose            | int.rate | installment | log.annual.inc | dti   | fico | d           | inq.last.6r | delinq.2yr | pub.rec | not.fully.paid |   |
|---------------|--------------------|----------|-------------|----------------|-------|------|-------------|-------------|------------|---------|----------------|---|
| 1             | debt_consolidation | 0.1189   | 829.1       | 11.35040654    | 19.48 | 737  | 5           | 0           | 0          | 0       | 0              |   |
| 2             | credit_card        | 0.1071   | 228.22      | 11.08214255    | 14.29 | 707  | 2           | 0           | 0          | 0       | 0              |   |
| 3             | debt_consolidation | 0.1357   | 366.86      | 10.37349118    | 11.63 | 682  | 2           | 0           | 0          | 0       | 0              |   |
| 4             | debt_consolidation | 0.1008   | 160.34      | 11.35040654    | 8.1   | 712  | 2           | 0           | 0          | 0       | 0              |   |
| 5             | credit_card        | 0.1436   | 102.92      | 11.29973234    | 14.97 | 667  | 2           | 0           | 0          | 0       | 0              |   |
| 6             | credit_card        | 0.0788   | 125.13      | 11.90486755    | 16.98 | 727  | 6           | 0           | 0          | 0       | 0              |   |
| 7             | debt_consolidation | 0.1496   | 194.02      | 10.71441777    | 4     | 667  | 3           | 0           | 0          | 1       | 1              |   |
| 8             | all_other          | 0.1114   | 131.22      | 11.00209984    | 11.08 | 722  | 2           | 0           | 0          | 0       | 1              |   |
| 9             | home_improvement   | 0.1134   | 87.19       | 11.40756495    | 17.25 | 682  | 2           | 1           | 0          | 0       | 0              |   |
| 10            | debt_consolidation | 0.1221   | 84.12       | 10.20359214    | 10    | 707  | 2           | 1           | 0          | 0       | 0              |   |
| 11            | debt_consolidation | 0.1347   | 360.43      | 10.4341158     | 22.09 | 677  | 6           | 2           | 0          | 1       | 0              |   |
| 12            | debt_consolidation | 0.1124   | 251.58      | 11.83500896    | 9.16  | 662  | 2           | 1           | 0          | 0       | 0              |   |
| 13            | debt_consolidation | 0.0859   | 316.11      | 10.93310697    | 15.49 | 767  | 6           | 0           | 0          | 0       | 0              |   |
| 14            | small_business     | 0.0714   | 92.82       | 11.51292546    | 6.5   | 747  | 2           | 0           | 1          | 0       | 0              |   |
| 15            | debt_consolidation | 0.0863   | 209.54      | 9.487972109    | 9.73  | 727  | 1559.958333 | 6282        | 44.6       | 0       | 0              | 0 |
| 16            | major_purchase     | 0.1103   | 327.53      | 10.73891524    | 13.04 | 702  | 8159.958333 | 5394        | 53.4       | 1       | 0              | 0 |
| 17            | all_other          | 0.1317   | 77.69       | 10.52277288    | 2.26  | 672  | 3895.958333 | 2211        | 88.4       | 0       | 0              | 0 |
| 18            | credit_card        | 0.0894   | 476.58      | 11.60823564    | 7.07  | 797  | 6510.958333 | 7586        | 52.7       | 1       | 0              | 0 |
| 19            | debt_consolidation | 0.1039   | 584.12      | 10.4921622     | 3.8   | 712  | 1553.958333 | 14354       | 58.6       | 0       | 2              | 0 |
| 20            | major_purchase     | 0.1513   | 173.65      | 11.00209984    | 2.74  | 667  | 1126.958333 | 591         | 84.4       | 3       | 0              | 0 |
| 21            | all_other          | 0.08     | 188.02      | 11.22524339    | 16.08 | 772  | 4888.958333 | 29797       | 23.2       | 1       | 0              | 0 |
| 22            | all_other          | 0.0863   | 474.42      | 10.81977828    | 2.59  | 797  | 11951       | 5656        | 27.6       | 0       | 0              | 0 |
| 23            | credit_card        | 0.1355   | 339.6       | 11.51292546    | 7.94  | 662  | 1939.958333 | 21162       | 57.7       | 0       | 0              | 0 |
| 24            | credit_card        | 0.0788   | 484.85      | 11.73606902    | 7.05  | 782  | 5640.041667 | 16931       | 34.6       | 1       | 0              | 0 |
| 25            | debt_consolidation | 0.1229   | 320.19      | 11.26444411    | 8.8   | 672  | 3760.958333 | 4822        | 58.1       | 0       | 1              | 0 |
| 26            | all_other          | 0.0901   | 159.03      | 12.4921622     | 10    | 712  | 1553.958333 | 14354       | 58.6       | 0       | 2              | 0 |
| 27            | all_other          | 0.0743   | 155.38      | 11.08214255    | 0.28  | 802  | 4649.958333 | 1576        | 5.7        | 1       | 0              | 0 |
| 28            | debt_consolidation | 0.1375   | 255.43      | 9.988797732    | 14.29 | 662  | 1318.958333 | 4175        | 51.5       | 0       | 1              | 0 |
| 29            | all_other          | 0.0743   | 155.38      | 12.20607265    | 0.28  | 772  | 4516.958333 | 3164        | 13.7       | 0       | 0              | 0 |
| 30            | all_other          | 0.0743   | 155.38      | 12.20607265    | 3.72  | 812  | 6778.958333 | 85607       | 0.7        | 0       | 0              | 0 |
| 31            | debt_consolidation | 0.0807   | 156.84      | 11.51292546    | 2.3   | 742  | 3148.958333 | 9698        | 19.4       | 0       | 0              | 0 |
| 32            | credit_card        | 0.1028   | 275.38      | 9.798127037    | 6.4   | 692  | 7469.958333 | 8847        | 26.9       | 1       | 1              | 0 |
| 33            | home_improvement   | 0.0743   | 155.38      | 11.91839057    | 0     | 777  | 7128.958333 | 6053        | 19.5       | 0       | 0              | 0 |
| 34            | home_improvement   | 0.0807   | 78.42       | 11.60823564    | 11.33 | 762  | 6559.958333 | 7274        | 13.1       | 0       | 0              | 0 |

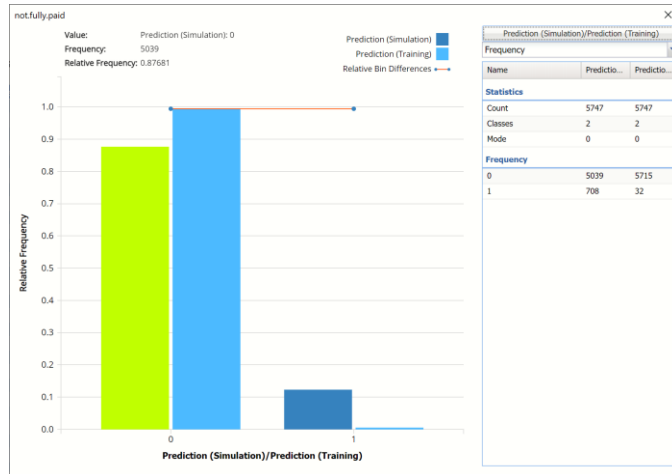
A very similar sequence of steps, selecting features of the dataset, partitioning the data into Training and Validation Sets as before, and choosing exactly the same parameters (sample size and the financial calculation) on the Simulation tab.

The screenshots show the following steps in the Analytic Solver Data Science interface:

- Logistic Regression - Parameters:** Shows the 'Data' tab with 'Partition Data' selected. The 'Partitioning' dialog box is open, showing 'Random partition' selected with a 'Set seed' of 12345. The 'Random partition percentages' are set to 80% for Training, 10% for Validation, and 10% for Test.
- Logistic Regression - Parameters:** Shows the 'Parameters' tab with 'Fit Intercept' checked and 'Iterations (Max)' set to 50.
- Logistic Regression - Scoring:** Shows the 'Scoring' tab with 'Score Training Data', 'Score Validation Data', and 'Score Test Data' all checked. The 'Summary Report' and 'Frequency Chart' options are also visible.
- Logistic Regression - Simulation:** Shows the 'Simulation' tab with 'Simulate Response Prediction' checked. The 'Fit Correlation' option is checked, and the 'Expression' field contains the formula:  $[[@installment]^12][@not.fully.paid]$ .

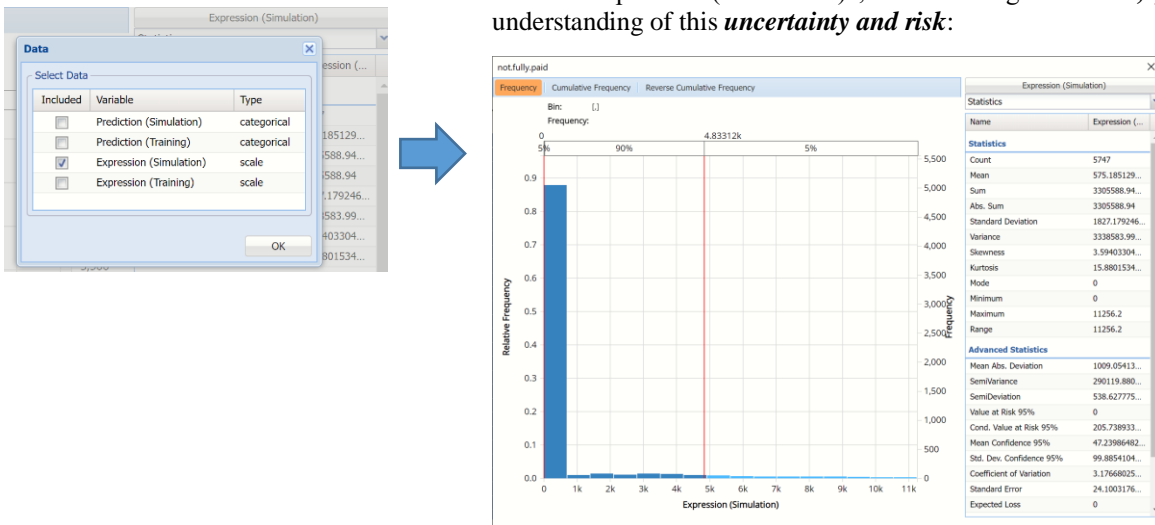
When Finish is clicked on the simulation tab, Analytic Solver Data Science again performs the “training” and “validation” steps to create a Logistic Regression ML model with parameters fitted to the data, and again renders the model in PMML (Predictive Modeling Markup Language) form.

This time, the chart on the Logreg\_Simulation tab in the Excel workbook shows quite different results:



Again, the light blue bars display the frequencies of “no default” (0) and “default” (1) on the 5,747 predicted cases in the Training Set. The dark blue bars display the predicted values in the simulated data: This shows that the Logistic Regression model’s future predictive power is quite good: It is predicting slightly more cases of “default” (1) and fewer cases of “no default”. Note the red Relative Bin Differences line that indicates that the absolute differences between the “no fault” and “default” bins is 0. Click the down arrow next to Frequency and select Bin Details for more information.

While the risk analysis of the Classification Tree model could “save” the user from an expensive error if he/she had put the Classification Tree model into production use, the much-better Logistic Regression model still has some uncertainty and risk. By clicking the button in the top right of the above chart, currently labeled “Prediction (Simulation) / Prediction (Training)”, and selecting instead “Expression (Simulation)”, the user can gain a *better, quantified* understanding of this *uncertainty and risk*:



This form of the chart – still a frequency histogram, but of the user-specified *financial consequence* – shows that in most cases the loan won’t default and the loss will be \$0, but *across all cases* the loss will average \$575 (the Mean value) and may be as large as \$11,256 (the Maximum value) in some cases. Other statistics such as Standard Deviation and Conditional Value at Risk give a quantitatively-oriented business manager some further insights into the model’s performance. The insights are available *before* the model is put into production

use, and can be useful in determining how and when the model's performance should be evaluated, and the model potentially re-trained in the future.



# Fitting the Best Model

---

## Introduction

Analytic Solver Data Science includes comprehensive, powerful support for data science and machine learning. Using these tools, you can "train" or fit your data to a wide range of statistical and machine learning models: Classification and regression trees, neural networks, linear and logistic regression, discriminant analysis, naïve Bayes, k-nearest neighbors and more. But the task of choosing and comparing these models, and selecting parameters for each one was up to you.

With the Find Best Model options, you can automate this work as well! Find Best Model uses methods similar to those in (expensive high-end) tools like DataRobot and RapidMiner, to automatically choose types of ML models and their parameters, validate and compare them according to criteria that you choose, and deliver the model that best fits your data.

Continue reading to discover how to utilize Find Best Model through an easy-to-follow example. For more information on how to run each classification or regression learner independently of Find Best Model, see the Data Science Reference Guide.

---

## Find Best Model Classification Example

This example illustrates how to utilize the Find Best Model method for Classification, included in Analytic Solver Data Science for Desktop Excel or Excel Online, by using the [Heart Failure Clinical Records Dataset](#)<sup>6</sup>. This dataset contains thirteen variables describing 299 patients experiencing heart failure. Find Best Model fits a model to all selected supervised learning (classification) methods in order to observe which method provides the best fit to the data. The goal of this example is to fit the best model to the dataset, then use this fitted model to determine if a new patient is at risk of perishing due to heart failure.

A description of each variable contained in the dataset appears in the table below.

| VARIABLE               | DESCRIPTION                                         |
|------------------------|-----------------------------------------------------|
| AGE                    | Age of patient                                      |
| ANAEMIA                | Decrease of red blood cells or hemoglobin (boolean) |
| CREATINE_PHOSPHOKINASE | Level of the CPK enzyme in the blood (mcg/L)        |
| DIABETES               | If the patient has diabetes (boolean)               |

---

<sup>6</sup> Davide Chicco, Giuseppe Jurman: Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone. BMC Medical Informatics and Decision Making 20, 16 (2020). ([link](#))

|                     |                                                                        |
|---------------------|------------------------------------------------------------------------|
| EJECTION_FRACTION   | Percentage of blood leaving the heart at each contraction (percentage) |
| HIGH_BLOOD_PRESSURE | If the patient has hypertension (boolean)                              |
| PLATELETS           | Platelets in the blood (kiloplatelets/mL)                              |
| SERUM_CREATININE    | Level of serum creatinine in the blood (mg/dL)                         |
| SERUM_SODIUM        | Level of serum sodium in the blood (mEq/L)                             |
| SEX                 | Woman (0) or man (1)                                                   |
| SMOKING             | If the patient smokes or not (boolean)                                 |
| TIME                | Follow-up period (days)                                                |
| DEATH_EVENT         | If the patient deceased during the follow-up period (boolean)          |

All supervised algorithms include a new Simulation tab. This tab uses the functionality from the Generate Data feature (described in the What's New section of this guide and then more in depth in the Analytic Solver Data Science Reference Guide) to generate synthetic data based on the training partition, and uses the fitted model to produce predictions for the synthetic data. The resulting report, CFBM\_Simulation, will contain the synthetic data, the predicted values and the Excel-calculated Expression column, if present. In addition, frequency charts containing the Predicted, Training, and Expression (if present) sources or a combination of any pair may be viewed, if the charts are of the same type. Since this new functionality does not support categorical variables, these types of variables will not be present in the model, only continuous variables.

## Opening the Dataset

Open the `Heart_failure_clinical_records_dataset.xlsx` by clicking **Help – Example Models – Forecasting/Data Science Examples**.

## Partitioning the Dataset

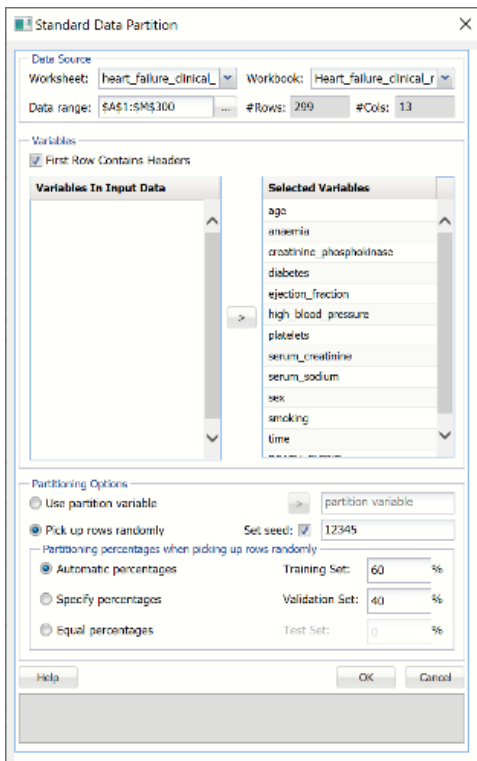
Partition the dataset by clicking **Partition – Standard Partition**.

- **Move all features under Variables In Input Data to Selected Variables.**
- Click **OK** to accept the partitioning defaults and create the partitions.

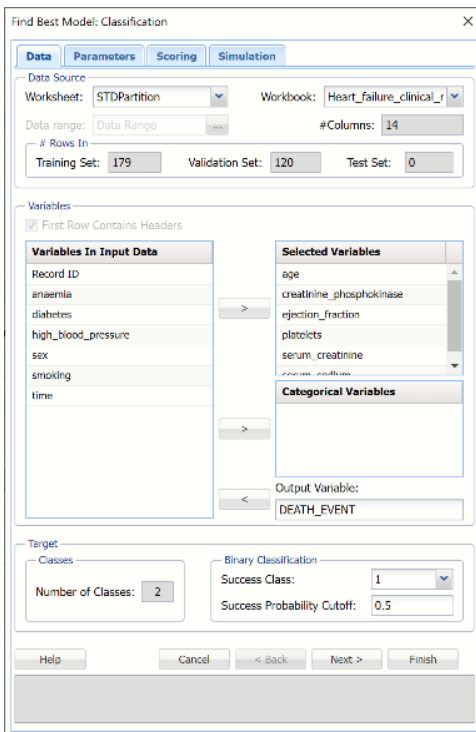
A new worksheet, `STDPartition`, is inserted directly to the right of the dataset. Find Best Model will be performed on both the Training and Validation partitions.

For more information on partitioning a dataset, see the Partitioning chapter within the Data Science Reference Guide.

Standard Data Partition Dialog



Find Best Model: Classification Data Tab



## Running Find Best Model

With STDPartition worksheet selected, click **Classify – Find Best Model** to open the Find Best Model Data tab.

### Data Tab

The continuous variables are selected on the Data tab, along with the Success Class and Success Probability Cutoff.

- **Select age, creatinine\_phosphokinase, ejection\_fraction, platelets, serum\_creatinine and serum\_sodium for Selected Variables.**
- **Select Death\_Event as Output Variable.**
- Leave Success Class and Success Probability Cutoff at their defaults.
- Click **Next** to move to the Parameters tab.

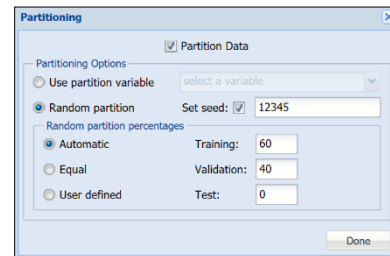
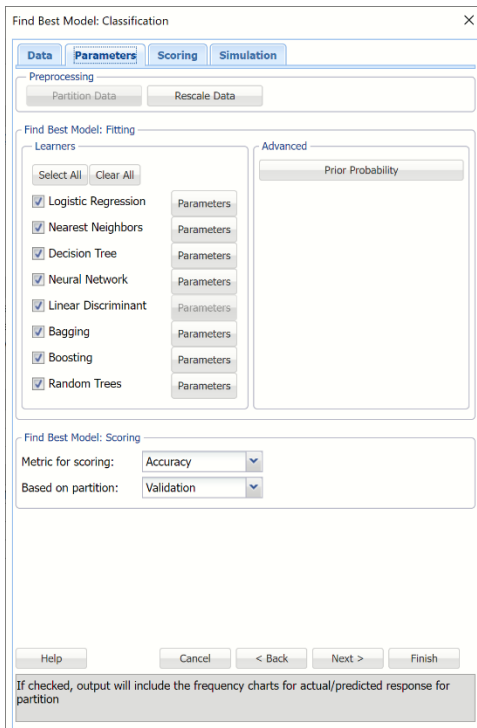
### Parameters tab

By default, all eligible supervised learners are automatically enabled based on the presence of categorical features or binary/multiclass classification. Optionally, all possible parameters for each algorithm may be defined using the Parameters button to the right of each learner.

### Partition Data

The Partition Data button is disabled because the original dataset has already been partitioned. Optionally, partitioning may be controlled from the Parameters tab, if desired.

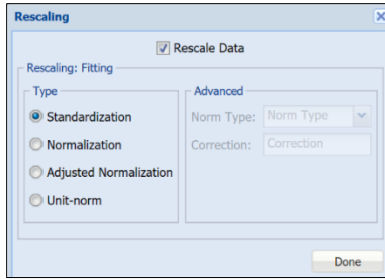
Find Best Model: Classification Parameters Tab



### Rescale Data

To rescale the data, click Rescale Data and select the Rescale Data option at the top of the Rescaling dialog.

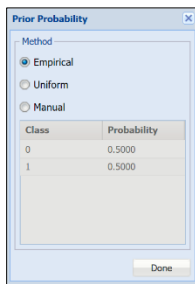
Use Rescaling to normalize one or more features in your data during the data preprocessing stage. Analytic Solver Data Science provides the following methods for feature scaling: Standardization, Normalization, Adjusted Normalization and Unit-norm. For more information on this feature, see the Rescale Continuous Data section within the Transform Continuous Data chapter that occurs in the Data Science Reference Guide.



- This example does not use rescaling to rescale the dataset data. Uncheck the Rescale Data option at the top of the dialog and click Done.

### Prior Probability

Click **Prior Probability**. Three options appear in the *Prior Probability* Dialog: *Empirical*, *Uniform* and *Manual*.



If the first option is selected, *Empirical*, Analytic Solver Data Science will assume that the probability of encountering a particular class in the dataset is the same as the frequency with which it occurs in the training data.

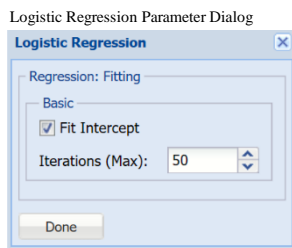
If the second option is selected, *Uniform*, Analytic Solver Data Science will assume that all classes occur with equal probability.

Select the third option, *Manual*, to manually enter the desired probability values for each class.

- Leave the default setting, *Empirical*, selected and click Done to close the dialog.

### Find Best Model: Fitting

To set a parameter for each selected learner, click the Parameters button to the right. The following chart gives a brief description of the parameters for each learner.

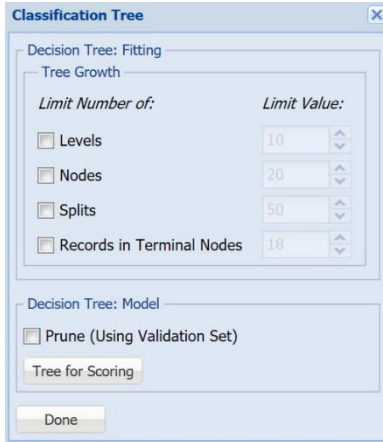


|                     |                                                                                                                                                                                                                              |
|---------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Logistic Regression | For more information on each parameter, see the Logistic Regression chapter within the Data Science Reference Guide.                                                                                                         |
| Fit Intercept       | When this option is selected, the default setting, Analytic Solver Data Science will fit the Logistic Regression intercept. If this option is not selected, Analytic Solver Data Science will force the intercept term to 0. |
| Iterations (Max)    | Estimating the coefficients in the Logistic Regression algorithm requires an iterative non-linear maximization procedure. You                                                                                                |

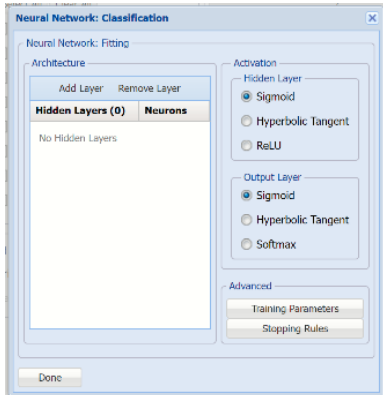
K-Nearest Neighbors Parameter Dialog



Classification Tree Parameter Dialog

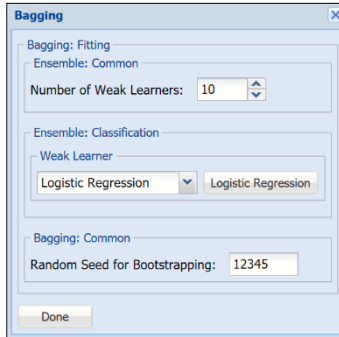


Neural Network Parameter Dialog

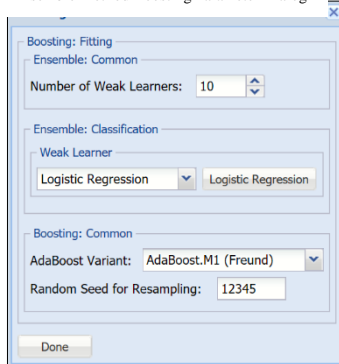


|                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                             |
|-------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                                                   | can specify a maximum number of iterations to prevent the program from getting lost in very lengthy iterative loops. This value must be an integer greater than 0 or less than or equal to 100 (1 < value <= 100).                                                                                                                                                                                                                          |
| K-Nearest Neighbors                                               | For more information on each parameter, see the K-Nearest Neighbors Classification chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                         |
| # Neighbors (k)                                                   | Enter a value for the parameter K in the Nearest Neighbor algorithm.                                                                                                                                                                                                                                                                                                                                                                        |
| Classification Tree                                               | For more information on each parameter, see the Classification Tree chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                        |
| Tree Growth Levels, Nodes, Splits, Tree Records in Terminal Nodes | In the <i>Tree Growth</i> section, select Levels, Nodes, Splits, and Records in Terminal Nodes. Values entered for these options limit tree growth, i.e. if 10 is entered for Levels, the tree will be limited to 10 levels.                                                                                                                                                                                                                |
| Prune                                                             | If a validation partition exists, this option is enabled. When this option is selected, Analytic Solver Data Science will prune the tree using the validation set. Pruning the tree using the validation set reduces the error from over-fitting the tree to the training data.<br><br>Click Tree for Scoring to click the Tree type used for scoring: Fully Grown, Best Pruned, Minimum Error, User Specified or Number of Decision Nodes. |
| Neural Network                                                    | For more information on each parameter, see the Neural Network Classification chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                              |
| Architecture                                                      | Click <i>Add Layer</i> to add a hidden layer. To delete a layer, click <i>Remove Layer</i> . Once the layer is added, enter the desired Neurons.                                                                                                                                                                                                                                                                                            |
| Hidden Layer                                                      | Nodes in the hidden layer receive input from the input layer. The output of the hidden nodes is a weighted sum of the input values. This weighted sum is computed with weights that are initially set at random values. As the network “learns”, these weights are adjusted. This weighted sum is used to compute the hidden node’s output using a <i>transfer function</i> . The default selection is <i>Sigmoid</i> .                     |
| Output Layer                                                      | As in the hidden layer output calculation (explained in the above paragraph), the output layer is also computed using the same transfer function as described for <i>Activation: Hidden Layer</i> . The default selection is <i>Sigmoid</i> .                                                                                                                                                                                               |
| Training Parameters                                               | Click Training Parameters to open the Training Parameters dialog to specify parameters related to the training of the Neural Network algorithm.                                                                                                                                                                                                                                                                                             |
| Stopping Rules                                                    | Click Stopping Rules to open the Stopping Rules dialog. Here users can specify a comprehensive set of rules for stopping the                                                                                                                                                                                                                                                                                                                |

Ensemble Method Bagging Parameter Dialog

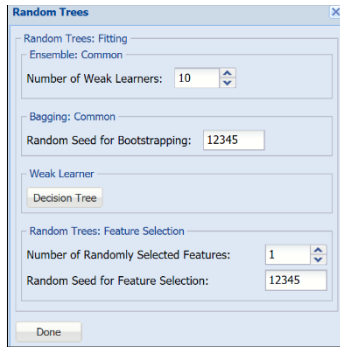


Ensemble Method Boosting Parameter Dialog



|                               |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
|-------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                               | algorithm early plus cross-validation on the training error.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Linear Discriminant           | For more information on this learner, see the Discriminant Analysis Classification chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| Bagging Ensemble Method       | For more information on each parameter, see the Ensemble Methods Classification chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
| Number of Weak Learners       | This option controls the number of “weak” classification models that will be created. The ensemble method will stop when the number of classification models created reaches the value set for this option. The algorithm will then compute the weighted sum of votes for each class and assign the “winning” classification to each record.                                                                                                                                                                                                                                                                                                                                                                                                             |
| Weak Learner                  | Under Ensemble: Classification click the down arrow beneath Weak Learner to select one of the six featured classifiers: Linear Discriminant, Logistic Regression, k-NN, Naïve Bayes, Neural Networks, or Decision Trees. The command button to the right will be enabled. Click this command button to control various option settings for the weak learner.                                                                                                                                                                                                                                                                                                                                                                                             |
| Random Seed for Bootstrapping | Enter an integer value to specify the seed for random resampling of the training data for each weak learner. Setting the random number seed to a nonzero value (any number of your choice is OK) ensures that the same sequence of random numbers is used each time the dataset is chosen for the classifier. The default value is “12345”. If left blank, the random number generator is initialized from the system clock, so the sequence of random numbers will be different in each calculation. If you need the results from successive runs of the algorithm to another to be strictly comparable, you should set the seed. To do this, type the desired number you want into the box. This option accepts positive integers with up to 9 digits. |
| Boosting Ensemble Method      | For more information on each parameter, see the Boosting Classification Ensemble Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Number of Weak Learners       | See description above.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
| Weak Learner                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
| Adaboost Variant              | In AdaBoost.M1 (Freund), the constant is calculated as:<br>$\alpha_b = \ln((1 - e_b) / e_b)$<br>In AdaBoost.M1 (Breiman), the constant is calculated as:<br>$\alpha_b = 1/2 \ln((1 - e_b) / e_b)$<br>In SAMME, the constant is calculated as:<br>$\alpha_b = 1/2 \ln((1 - e_b) / e_b + \ln(k - 1))$ where k is the number of classes                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Random Seed for Resampling    | Enter an integer value to specify the seed for random resampling of the training data for                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |

Ensemble Method Random Trees Parameter Dialog



|                                                                                  |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
|----------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                                                                  | each weak learner. Setting the random number seed to a nonzero value (any number of your choice is OK) ensures that the same sequence of random numbers is used each time the dataset is chosen for the classifier. The default value is “12345”.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Random Trees Ensemble Method                                                     | For more information on each parameter, see the Boosting Classification Ensemble Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Number of Weak Learners<br><br>Random Seed for Bootstrapping<br><br>Weak Learner | See description above.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Number of Randomly Selected Features                                             | The Random Trees ensemble method works by training multiple “weak” classification trees using a fixed number of randomly selected features then taking the mode of each class to create a “strong” classifier. The option <i>Number of randomly selected features</i> controls the fixed number of randomly selected features in the algorithm. The default setting is <b>3</b> .                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Feature Selection Random Seed                                                    | If an integer value appears for <i>Feature Selection Random seed</i> , Analytic Solver Data Science will use this value to set the feature selection random number seed. Setting the random number seed to a nonzero value (any number of your choice is OK) ensures that the same sequence of random numbers is used each time the dataset is chosen for the classifier. The default value is “12345”. If left blank, the random number generator is initialized from the system clock, so the sequence of random numbers will be different in each calculation. If you need the results from successive runs of the algorithm to another to be strictly comparable, you should set the seed. To do this, type the desired number you want into the box. This option accepts positive integers with up to 9 digits. |

### Find Best Model: Scoring

The two parameters at the bottom of the dialog under Find Best Model: Scoring, determine how well each classification method fits the data.

The Metric for Scoring may be changed to Accuracy, Specificity, Sensitivity, Precision or F-1. See the table below for a brief description of each statistic.

- For this example, leave Accuracy selected.

| Statistic   | Description                                                                                       |
|-------------|---------------------------------------------------------------------------------------------------|
| Accuracy    | Accuracy refers to the ability of the classifier to predict a class label correctly.              |
| Specificity | Specificity is defined as the proportion of negative classifications that were actually negative. |

|             |                                                                                                                                                                                                                                                                                                                                                                 |
|-------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Sensitivity | Sensitivity is defined as the proportion of positive cases there were classified correctly as positive.                                                                                                                                                                                                                                                         |
| Precision   | Precision is defined as the proportion of positive results that are truly positive.                                                                                                                                                                                                                                                                             |
| F-1         | <p>Calculated as <math>2 \times (\text{Precision} * \text{Sensitivity}) / (\text{Precision} + \text{Sensitivity})</math></p> <p>The F-1 Score provides a statistic to balance between Precision and Sensitivity, especially if an uneven class distribution exists. The closer the F-1 score is to 1 (the upper bound) the better the precision and recall.</p> |

The options for Based on Partition depend on the number of partitions present. In this example, the original dataset was partitioned into training and validation partitions so those will be the options present in the drop down menu.

- For this example **select Validation**.

Click Next to advance to the Scoring tab.

### Scoring tab

Output options are selected on the Scoring tab. In this example select all four options: Detailed Report, Summary Report and Lift Charts and Frequency Chart under both Score Training Data and Score Validation Data.

By default, CFBM\_Output and CFBM\_Stored are generated and inserted directly to the right of the STDPartition.

- CFBM\_Output contains a listing of all model inputs such a input/output variables and parameter settings for all Learners, as well as Model Performance tables containing evaluations for every available metric, every learner on all available partitions. The Learner identified internally as performing the best, is highlighted in red. (Recall that the statistic used for determining which Learner performs best on the dataset was selected on the Parameters tab.)
- FBM\_Stored contains the PMML (Predictive Model Markup Language) model which can be utilized to score new data. For more information on scoring, see the Scoring chapter that appears later in this guide.

Selecting *Detailed Report* produces CFBM\_TrainingScore and CFBM\_ValidationScore.

- Both reports contain detailed scoring information on both the Training and Validation partitions using the "best" learner.

Misclassified records are highlighted in red.

*Summary Report* is selected by default. This option produces a summary report at the top of both CFBM\_TrainingScore and CFBM\_ValidationScore worksheets.

- *Summary Reports* contains a Confusion Matrix, Error Report and various metrics. The chart to the left displays how each metric is calculated.



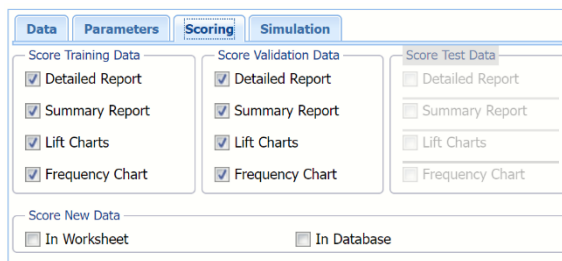
- True Positive is the number of real positives in the dataset (i.e. the number of positive records that were classified as positive by the learner).
- True Negative is the number of real negatives in the dataset (i.e. the number of negative records that were classified as negative by the learner).
- False Negative is the number of records classified as negative that are actually positive in the dataset.
- False Positive is the number of records classified as positive that are actually negative in the dataset.
- Selecting *Lift Charts* generates Lift Charts, ROC Curves and Decile-Wise Lift Charts.
- When Frequency Chart is selected, a frequency chart will be displayed when the CFBM\_TrainingScore and CFBM\_ValidationScore worksheets are selected. This chart will display an interactive application similar to the Analyze Data feature, explained in detail in the Analyze Data chapter that appears earlier in this guide. This chart will include frequency distributions of the actual and predicted responses individually, or side-by-side, depending on the user’s preference, as well as basic and advanced statistics for variables, percentiles, six sigma indices.

Since we did not create a test partition, the options for Score test data are disabled. See the chapter “Data Science Partitioning” for information on how to create a test partition.

See the *Scoring New Data* chapter that appears later in this guide for more information on *Score New Data* in options.

Click Finish to run Find Best Model.

Find Best Model Classification Scoring tab

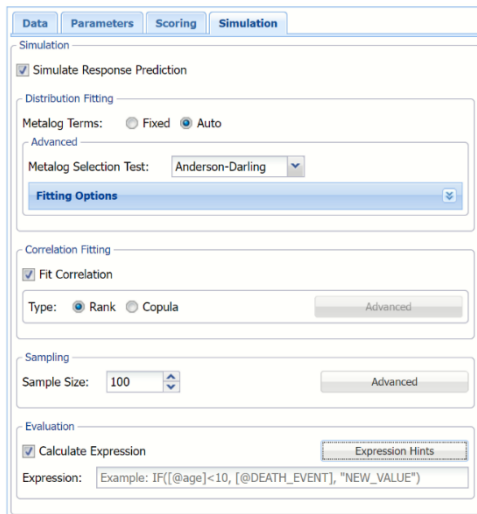


Click **Next** to advance to the Simulation tab.

Select Simulation Response Prediction to enable all options on the Simulation tab of the Find Best Model Classification dialog.

**Simulation tab:** All supervised algorithms include a new Simulation tab. This tab uses the functionality from the Generate Data feature (described earlier in this guide) to generate synthetic data based on the training partition, and uses the fitted model to produce predictions for the synthetic data. The resulting report, CFMB\_Simulation, will contain the synthetic data, the predicted values and the Excel-calculated Expression column, if present. In addition, frequency charts containing the Predicted, Training, and Expression (if present) sources or a combination of any pair may be viewed, if the charts are of the same type.

*Find Best Model Classification dialog, Simulation tab*



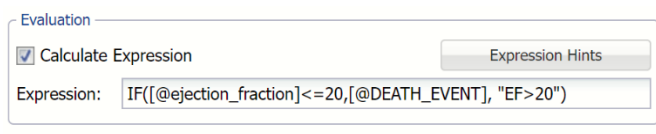
**Evaluation:** Select Calculate Expression to amend an Expression column onto the frequency chart displayed on the CFBM\_Simulation output tab. Expression can be any valid Excel formula that references a variable and the response as [COLUMN\_NAME]. Click the Expression Hints button for more information on entering an expression.

For the purposes of this example, leave all options at their defaults in the Distribution Fitting, Correlation Fitting and Sampling sections of the dialog. For Expression, enter the following formula to display if the patient suffered catastrophic heart failure (@DEATH\_EVENT) when his/her Ejection\_Fraction was less than or equal to 20.

IF([@ejection\_fraction]<=20,[@DEATH\_EVENT], "EF>20")

Note that variable names are case sensitive.

*Evaluation section on the Find Best Model dialog, Simulation tab*



For more information on the remaining options shown on this dialog in the Distribution Fitting, Correlation Fitting and Sampling sections, see the Generate Data chapter that appears earlier in this guide.

Click **Finish** to run Find Best Model on the example dataset.

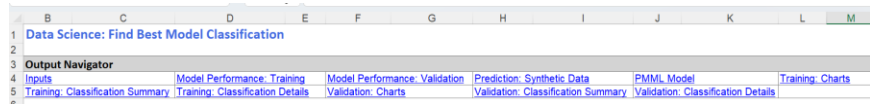
## Interpreting the Results

All worksheets containing results from Find Best Model are inserted to the right of the STDPartition worksheet.

### CFBM\_Output Result Worksheet

The CFBM\_Output worksheet is inserted directly to the right of the STDPartition worksheet. This report lists all input variables and all parameter settings for each learner, along with the Model Performance of each Learner on all partitions. This example utilizes two partitions, training and validation.

The Output Navigator appears at the very top of this worksheet. Click the links to easily move between each section of the output. The Output Navigator is listed at the top of each worksheet included in the output.



The Inputs section includes information pertaining to the dataset, the input variables and parameter settings.

| Data                                      |                                             |
|-------------------------------------------|---------------------------------------------|
| Worksheet                                 | Heart_failure_clinical_records_dataset.xlsx |
| Training data used for building the model | SC537:SP5215                                |
| # Records in the training data            | 179                                         |
| Validation data                           | SC5216:SP5335                               |
| # Records in the validation data          | 120                                         |

| Variables             |                                                                    |
|-----------------------|--------------------------------------------------------------------|
| # Variables           | 6                                                                  |
| Scale Variables       | age creatinine_pl/ejection_fractiplatelets serum_creatserum_sodium |
| Categorical Variables |                                                                    |
| Output Variable       | DEATH_EVENT                                                        |

| Rescaling: Fitting Parameters |       |
|-------------------------------|-------|
| Rescale Data?                 | FALSE |

| Find Best Model: Fitting Parameters |                                             |
|-------------------------------------|---------------------------------------------|
| Learners                            | LOGISTIC_REGRESSION, NEAREST_NEIGHBORS, DEC |

| Find Best Model Classification: Fitting Parameters |           |
|----------------------------------------------------|-----------|
| Prior Probability Calculation                      | EMPIRICAL |

| Find Best Model Classification: Model Parameters |            |
|--------------------------------------------------|------------|
| Metric for scoring                               | Accuracy   |
| Based on partition                               | Validation |
| # Classes                                        | 2          |
| Success Class                                    | 1          |
| Success Probability                              | 0.5        |

Further down within Inputs, the parameter selections for each Learner are listed.

| Data                                      |                                             |
|-------------------------------------------|---------------------------------------------|
| Worksheet                                 | Heart_failure_clinical_records_dataset.xlsx |
| Training data used for building the model | SC537:SP5215                                |
| # Records in the training data            | 179                                         |
| Validation data                           | SC5216:SP5335                               |
| # Records in the validation data          | 120                                         |

| Variables             |                                                                    |
|-----------------------|--------------------------------------------------------------------|
| # Variables           | 6                                                                  |
| Scale Variables       | age creatinine_pl/ejection_fractiplatelets serum_creatserum_sodium |
| Categorical Variables |                                                                    |
| Output Variable       | DEATH_EVENT                                                        |

| Rescaling: Fitting Parameters |       |
|-------------------------------|-------|
| Rescale Data?                 | FALSE |

| Find Best Model: Fitting Parameters |                                             |
|-------------------------------------|---------------------------------------------|
| Learners                            | LOGISTIC_REGRESSION, NEAREST_NEIGHBORS, DEC |

| Find Best Model Classification: Fitting Parameters |           |
|----------------------------------------------------|-----------|
| Prior Probability Calculation                      | EMPIRICAL |

| Find Best Model Classification: Model Parameters |          |
|--------------------------------------------------|----------|
| Metric for scoring                               | Accuracy |
| Based on partition                               | Training |
| # Classes                                        | 2        |
| Success Class                                    | 1        |
| Success Probability                              | 0.5      |

| Logistic Regression                            |      |
|------------------------------------------------|------|
| <b>Regression Model: Fitting Parameters</b>    |      |
| Fit Intercept                                  | TRUE |
| <b>Logistic Regression: Fitting Parameters</b> |      |
| Maximum Number of Iterations                   | 150  |

| Nearest Neighbors                            |   |
|----------------------------------------------|---|
| <b>Nearest Neighbors: Fitting Parameters</b> |   |
| # Nearest neighbors (K)                      | 1 |

| Decision Tree                          |             |
|----------------------------------------|-------------|
| <b>Decision Tree: Model Parameters</b> |             |
| Prune?                                 | FALSE       |
| Scoring tree type                      | Fully grown |

| Neural Network                            |                  |
|-------------------------------------------|------------------|
| <b>Neural Network: Fitting Parameters</b> |                  |
| Random seed for initial weights           | 12345            |
| # Hidden Layers                           | 0                |
| Learning rate                             | 0.1              |
| Weight change momentum                    | 0.8              |
| Error tolerance                           | 0.01             |
| Weight decay                              | 0                |
| Cost function                             | Cross Entropy    |
| Hidden layer activation function          | LOGISTIC_Sigmoid |
| Output layer activation function          | LOGISTIC_Sigmoid |
| Learning order                            | Original         |
| Response cost function                    | 0.01             |
| Data for error computation                | TRAINING ONLY    |
| Maximum number of epochs                  | 30               |
| Maximum number of epochs without imp      | 5                |
| Maximum training time                     | 3600             |
| Minimum relative change in error          | 0.0001           |
| Minimum relative change in error comp     | 0.001            |

| Linear Discriminant                                |                     |
|----------------------------------------------------|---------------------|
| <b>Ensemble Parameters</b>                         |                     |
| Weak learner                                       | Classification Tree |
| Number of weak learners                            | 10                  |
| Show weak learner models?                          | FALSE               |
| <b>Bagging Parameters</b>                          |                     |
| Bootstrap seed                                     | 12345               |
| <b>Boosting</b>                                    |                     |
| <b>Ensemble Parameters</b>                         |                     |
| Weak learner                                       | Classification Tree |
| Number of weak learners                            | 10                  |
| Show weak learner models?                          | FALSE               |
| <b>Boosting Classification: Fitting Parameters</b> |                     |
| Adaptivity Variant                                 | ML-FREUND           |
| Resampling Seed                                    | 12345               |
| <b>Random Trees</b>                                |                     |
| <b>Ensemble Parameters</b>                         |                     |
| Weak learner                                       | Classification Tree |
| Number of weak learners                            | 10                  |
| Show weak learner models?                          | FALSE               |
| <b>Bagging Parameters</b>                          |                     |
| Bootstrap seed                                     | 12345               |
| <b>Random Trees: Fitting Parameters</b>            |                     |
| # Selected features                                | 2                   |
| Feature selection random seed                      | 12345               |

|    | C                                                  | D                                                 | E | F | G | H | I |
|----|----------------------------------------------------|---------------------------------------------------|---|---|---|---|---|
| 11 |                                                    |                                                   |   |   |   |   |   |
| 12 | <b>Simulation: Distribution Fitting Parameters</b> |                                                   |   |   |   |   |   |
| 13 | Metalog Terms                                      | Auto                                              |   |   |   |   |   |
| 14 | GOF Test                                           | Anderson-Darling                                  |   |   |   |   |   |
| 15 | Options                                            | {"age":{"lb":40,"numTerms":5,"ub":95},"creatinine |   |   |   |   |   |
| 16 |                                                    |                                                   |   |   |   |   |   |
| 17 | <b>Simulation: Correlation Fitting Parameters</b>  |                                                   |   |   |   |   |   |
| 18 | Correlation Type                                   | Rank                                              |   |   |   |   |   |
| 19 |                                                    |                                                   |   |   |   |   |   |
| 20 | <b>Simulation: Sampling Parameters</b>             |                                                   |   |   |   |   |   |
| 21 | Generate sample                                    | Yes                                               |   |   |   |   |   |
| 22 | Sample size                                        | 100                                               |   |   |   |   |   |
| 23 | Random seed                                        | 12345                                             |   |   |   |   |   |
| 24 | Random generator                                   | Mersenne Twister                                  |   |   |   |   |   |
| 25 | Sampling method                                    | Latin Hypercube                                   |   |   |   |   |   |
| 26 | Random streams                                     | Independent                                       |   |   |   |   |   |
| 27 | Calculate expression?                              | Yes                                               |   |   |   |   |   |
| 28 | Expression                                         | IF(@ejection_fraction<=20,@DEATH_EVENT),"EF       |   |   |   |   |   |
| 29 |                                                    |                                                   |   |   |   |   |   |

Scroll down to view any generated messages from the Find Best Model feature and also to view the performance of each learner on the training and validation partitions.

|     | B                                               | C | D | E | F | G |
|-----|-------------------------------------------------|---|---|---|---|---|
| 130 |                                                 |   |   |   |   |   |
| 137 | <b>Output Options</b>                           |   |   |   |   |   |
| 138 | Summary report of scoring on training data      |   |   |   |   |   |
| 139 | Detailed report of scoring on training data     |   |   |   |   |   |
| 140 | Lift charts on training data                    |   |   |   |   |   |
| 141 | Frequency chart on training data                |   |   |   |   |   |
| 142 | Summary report of scoring on validation data    |   |   |   |   |   |
| 143 | Detailed report of scoring on validation data   |   |   |   |   |   |
| 144 | Lift charts on validation data                  |   |   |   |   |   |
| 145 | Frequency chart on validation data              |   |   |   |   |   |
| 146 |                                                 |   |   |   |   |   |
| 147 |                                                 |   |   |   |   |   |
| 148 | <b>Note:</b> Scoring will be done using Bagging |   |   |   |   |   |

Metrics/Partition for Selecting Best Model

Random Trees Parameters

---

**Find Best Model: Scoring**

Metric for scoring: Accuracy

Based on partition: Validation

The Messages portion of the report indicates that Scoring will be performed using the Bagging ensemble method, the Learner selected as the "best" choice according to the selection for Find Best Model: Scoring parameters on the Parameters tab: Validation Partition Accuracy Metric.

Since the Bagging Accuracy metric for the Validation Partition has the highest score, that is the Learner that will be used for scoring.

|     | B                                    | C                   | D                   | E           | F           | G         | H         | I |
|-----|--------------------------------------|---------------------|---------------------|-------------|-------------|-----------|-----------|---|
| 150 | <b>Model Performance: Training</b>   |                     |                     |             |             |           |           |   |
| 151 |                                      |                     |                     |             |             |           |           |   |
| 152 | Metric                               | Accuracy (#correct) | Accuracy (%correct) | Specificity | Sensitivity | Precision | F1 score  |   |
| 153 | Logistic Regression                  | 137                 | 76.53631285         | 0.90598291  | 0.5         | 0.7380952 | 0.596154  |   |
| 154 | Decision Tree                        | 179                 | 100                 | 1           | 1           | 1         | 1         |   |
| 155 | Nearest Neighbors                    | 179                 | 100                 | 1           | 1           | 1         | 1         |   |
| 156 | Neural Network                       | 117                 | 65.36312849         | 1           | 0           | N/A       | N/A       |   |
| 157 | Linear Discriminant                  | 137                 | 76.53631285         | 0.90598291  | 0.5         | 0.7380952 | 0.596154  |   |
| 158 | Bagging                              | 173                 | 96.64804469         | 0.94871795  | 1           | 0.9117647 | 0.953846  |   |
| 159 | Boosting                             | 179                 | 100                 | 1           | 1           | 1         | 1         |   |
| 160 | Random Trees                         | 175                 | 97.76536313         | 0.97435897  | 0.98387097  | 0.953125  | 0.968254  |   |
| 161 |                                      |                     |                     |             |             |           |           |   |
| 162 | <b>Model Performance: Validation</b> |                     |                     |             |             |           |           |   |
| 163 |                                      |                     |                     |             |             |           |           |   |
| 164 | Metric                               | Accuracy (#correct) | Accuracy (%correct) | Specificity | Sensitivity | Precision | F1 score  |   |
| 165 | Logistic Regression                  | 91                  | 75.83333333         | 0.88372093  | 0.44117647  | 0.6       | 0.508475  |   |
| 166 | Decision Tree                        | 76                  | 63.33333333         | 0.63953488  | 0.61764706  | 0.4038462 | 0.488372  |   |
| 167 | Nearest Neighbors                    | 71                  | 59.16666667         | 0.68604651  | 0.35294118  | 0.3076923 | 0.328767  |   |
| 168 | Neural Network                       | 86                  | 71.66666667         | 1           | 0           | N/A       | N/A       |   |
| 169 | Linear Discriminant                  | 88                  | 73.33333333         | 0.88372093  | 0.35294118  | 0.5454545 | 0.428571  |   |
| 170 | Bagging                              | 94                  | 78.33333333         | 0.79069767  | 0.76470588  | 0.5909091 | 0.6666667 |   |
| 171 | Boosting                             | 84                  | 70                  | 0.75581395  | 0.55882353  | 0.475     | 0.513514  |   |
| 172 | Random Trees                         | 83                  | 69.16666667         | 0.73255814  | 0.58823529  | 0.4651163 | 0.519481  |   |
| 173 |                                      |                     |                     |             |             |           |           |   |

## CFBM\_TrainingScore

CFBM\_TrainingScore contains the Classification Summary and the Classification Details reports. Both reports have been generated using the Bagging Learner, as discussed above.

## CFBM\_TrainingScore

Click the *CFBM\_TrainingScore* tab to view the newly added Output Variable frequency chart, the Training: Classification Summary and the Training:

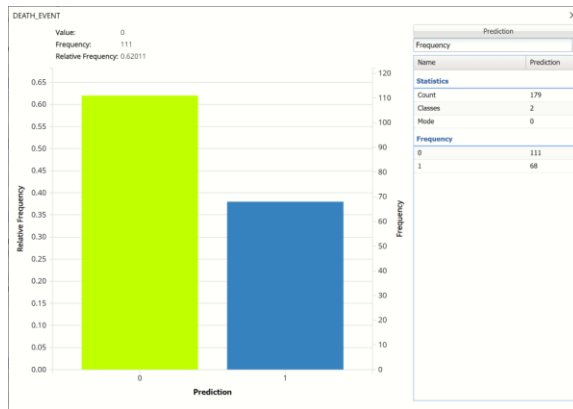
Classification Details report. All calculations, charts and predictions on this worksheet apply to the Training data.

Note: To view charts in the Cloud app, click the Charts icon on the Ribbon, select a worksheet under Worksheet and a chart under Chart.

- Frequency Charts:** The output variable frequency chart opens automatically once the *CFBM\_TrainingScore* worksheet is selected. To close this chart, click the “x” in the upper right hand corner of the chart. To reopen, click onto another tab and then click back to the *CFBM\_TrainingScore* tab. To move the chart, grab the dialog’s title bar and drag the chart to the desired location on the screen.

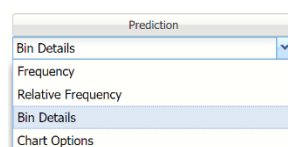
**Frequency:** This chart shows the frequency for both the predicted and actual values of the output variable, along with various statistics such as count, number of classes and the mode.

*Frequency Chart on CFBM\_TrainingScore output sheet*



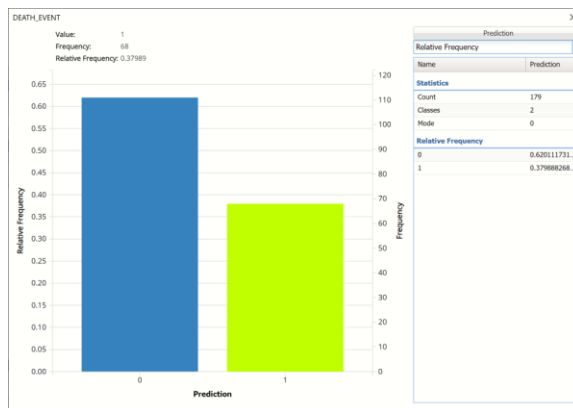
Click the down arrow next to Frequency to switch to Relative Frequency, Bin Details or Chart Options view.

*Frequency Chart, Frequency View*



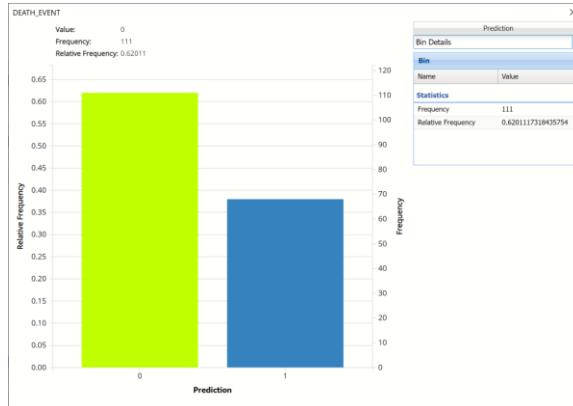
**Relative Frequency:** Displays the relative frequency chart.

*Relative Frequency Chart*



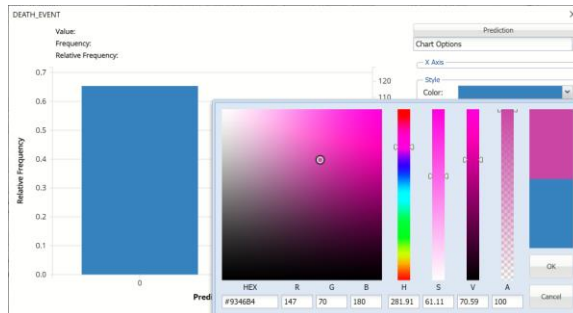
**Bin Details:** Use this view to find metrics related to each bin in the chart.

*Bin Details Chart*



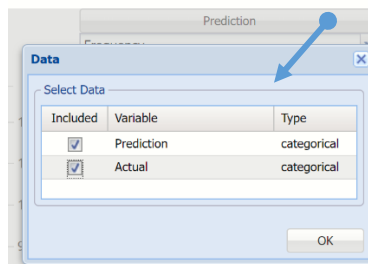
**Chart Options:** Use this view to change the color of the bars in the chart.

*Chart Options View*

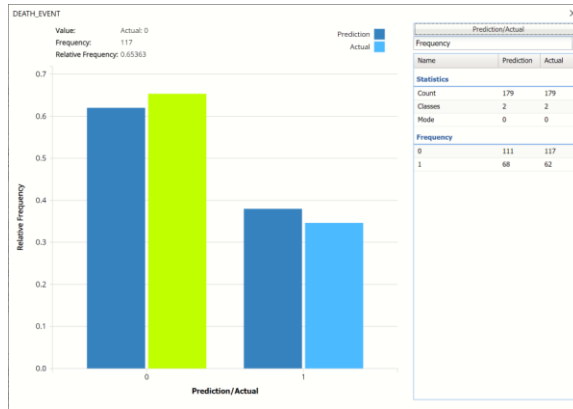


To see both predicted and actual values for the output variable in the training partition, click Prediction and select Actual. This change will be reflected on all charts.

*Click Prediction to change view*



*Frequency Chart displaying both Prediction and Actual data*



- **Classification Summary:** In the Classification Summary report, a Confusion Matrix is used to evaluate the performance of the classification method.

| Confusion Matrix |                 |    |
|------------------|-----------------|----|
|                  | Predicted Class |    |
| Actual Class     | 1               | 0  |
| 1                | TP              | FN |
| 0                | FP              | TN |

- TP stands for True Positive. These are the number of cases classified as belonging to the Success class that actually were members of the Success class.
- FN stands for False Negative. These are the number of cases that were classified as belonging to the Failure class when they were actually members of the Success class
- FP stands for False Positive. These cases were assigned to the Success class but were actually members of the Failure group
- TN stands for True Negative. These cases were correctly assigned to the Failure group.

*Confusion Matrix for Training Partition*

| Training: Classification Summary |             |          |             |
|----------------------------------|-------------|----------|-------------|
| <b>Confusion Matrix</b>          |             |          |             |
| Actual\Predicted                 | 0           | 1        |             |
| 0                                | 111         | 6        |             |
| 1                                | 0           | 62       |             |
| <b>Error Report</b>              |             |          |             |
| Class                            | # Cases     | # Errors | % Error     |
| 0                                | 117         | 6        | 5.128205128 |
| 1                                | 62          | 0        | 0           |
| Overall                          | 179         | 6        | 3.351955307 |
| <b>Metrics</b>                   |             |          |             |
| Metric                           | Value       |          |             |
| Accuracy (#correct)              | 173         |          |             |
| Accuracy (%correct)              | 96.64804469 |          |             |
| Specificity                      | 0.948717949 |          |             |
| Sensitivity (Recall)             | 1           |          |             |
| Precision                        | 0.911764706 |          |             |
| F1 score                         | 0.953846154 |          |             |
| Success Class                    | 1           |          |             |
| Success Probability              | 0.5         |          |             |

The Confusion Matrix for the Training Partition indicates:

- True Negative = 111 – These records are truly negative records. All 111 patients are survivors
- True Positive = 62 – These records are truly positive. All 62 patients are deceased.
- False Positive = 6 – These negative records were misclassified as positive. In other words, 6 surviving patients were misclassified as non-survivors.
- False Negative = 0 – No positive records were misclassified as negative, or there were no deceased patients that were misclassified as survivors.

The %Error for the training partition was 3.35%. All of the error is due to False Positives. The Bagging Ensemble misclassified 3.35% (6/179) of surviving patients as non-survivors.

The Metrics indicate:

- Accuracy: The number of records classified correctly by the Bagging Ensemble is 173 (111 + 62).
- %Accuracy: The percentage of records classified correctly is 96.1% (173/179).
- Specificity:  $(\text{True Negative})/(\text{True Negative} + \text{False Positives})$   
 $111/(111 + 6) = 0.949$

Specificity is defined as the proportion of negative classifications that were actually negative, or the fraction of survivors that actually survived. In this model, 111 actual surviving patients were classified correctly as survivors. There were 6 false positives or 6 actual surviving patients that were classified as deceased.

- Sensitivity or Recall:  $(\text{True Positive})/(\text{True Positive} + \text{False Negative})$   
 $62/(62 + 0) = 1.0$

Sensitivity is defined as the proportion of positive cases there were classified correctly as positive, or the proportion of actually deceased patients there were classified as deceased. In this model, 62 actual deceased patients were correctly classified as deceased. There were no false negatives or no actual deceased patients were incorrectly classified as a survivor.

Note: Since the object of this model is to correctly classify which patients will succumb to heart failure, this is an important statistic as it is very important for a physician to be able to accurately predict which patients require mitigation.

- Precision:  $(\text{True Positives})/(\text{True Positives} + \text{False Positives})$   
 $62/(62 + 6) = 0.912$

Precision is defined as the proportion of positive results that are true positive. In this model, 62 actual deceased patients were classified correctly as deceased. There were 6 false positives or 6 actual survivors classified incorrectly as deceased.

- F-1 Score:  $2 \times (\text{Precision} * \text{Sensitivity})/(\text{Precision} + \text{Sensitivity})$   
 $2 \times (0.912 * 0.949) / (0.912 + 0.949) = 0.954$



The F-1 Score provides a statistic to balance between Precision and Sensitivity, especially if an uneven class distribution exists, as in this example, (117 survivors vs 62 deceased). The closer the F-1 score is to 1 (the upper bound) the better the precision and recall.

- Success Class and Success Probability simply reports the settings for these two values as input on the Find Best Model: Classification Data tab.
- **Classification Details:** Individual records and their classifications are shown beneath Training: Classification Details. Any misclassified records are highlighted in red.

| Record ID | DEATH_EVENT | Prediction: DEATH_EVENT | PostProb: 0 | PostProb: 1 |
|-----------|-------------|-------------------------|-------------|-------------|
| Record 1  | 1           | 1                       | 0           | 1           |
| Record 5  | 1           | 1                       | 0           | 1           |
| Record 8  | 1           | 1                       | 0.2         | 0.8         |
| Record 15 | 0           | 0                       | 0.6         | 0.4         |
| Record 18 | 1           | 1                       | 0.1         | 0.9         |
| Record 20 | 1           | 1                       | 0.2         | 0.7         |

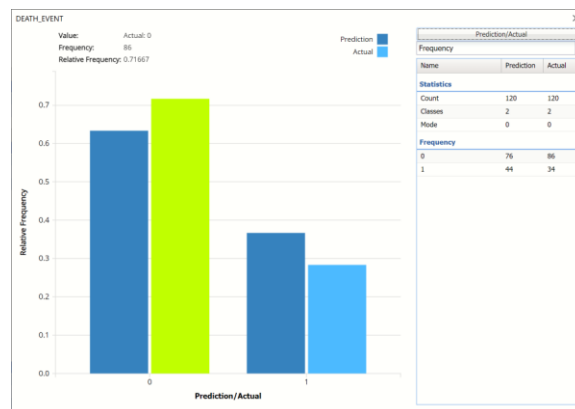
### CFBM\_ValidationScore

Click the *CFBM\_ValidationScore* tab to view the newly added Output Variable frequency chart, the Validation: Classification Summary and the Validation: Classification Details report. All calculations, charts and predictions on this worksheet apply to the Validation data.

- **Frequency Charts:** The output variable frequency chart opens automatically once the CFBM\_ValidationScore worksheet is selected. To close this chart, click the “x” in the upper right hand corner. To reopen, click onto another tab and then click back to the CFBM\_ValidationScore tab.

As explained above, this chart displays the frequency for both the predicted and actual values of the output variable, along with various statistics such as count, number of classes and the mode. Select Relative Frequency from the drop down menu, on the right, to see the relative frequencies of the output variable for both actual and predicted. See above for more information on this chart.

*CFBM\_ValidationScore Frequency Chart*



The Confusion Matrix for the Validation Partition indicates:

- True Negative = 68 – These records are truly negative records. All 68 patients are survivors
- True Positive = 26 – These records are truly positive. All 26 patients are deceased.
- False Positive = 18 – These negative records were misclassified as positive. In other words, 18 surviving patients were misclassified as non-survivors.
- False Negative = 8 – These positive records were misclassified as negative, or these 8 deceased patients were misclassified as survivors.

*Validation Partition Classification Summary*

| Validation: Classification Summary |             |          |             |
|------------------------------------|-------------|----------|-------------|
| <b>Confusion Matrix</b>            |             |          |             |
| Actual\Predicted                   | 0           | 1        |             |
| 0                                  | 68          | 18       |             |
| 1                                  |             | 8        | 26          |
| <b>Error Report</b>                |             |          |             |
| Class                              | # Cases     | # Errors | % Error     |
| 0                                  | 86          | 18       | 20.93023256 |
| 1                                  | 34          | 8        | 23.52941176 |
| Overall                            | 120         | 26       | 21.66666667 |
| <b>Metrics</b>                     |             |          |             |
| Metric                             | Value       |          |             |
| Accuracy (#correct)                | 94          |          |             |
| Accuracy (%correct)                | 78.33333333 |          |             |
| Specificity                        | 0.790697674 |          |             |
| Sensitivity (Recall)               | 0.764705882 |          |             |
| Precision                          | 0.590909091 |          |             |
| F1 score                           | 0.666666667 |          |             |
| Success Class                      | 1           |          |             |
| Success Probability                | 0.5         |          |             |

The %Error for the validation partition was 21.67%. Most of the error is due to False Positives. The Bagging Ensemble misclassified 15% (18/120) of surviving patients as non-survivors and 6.7% of non-surviving patients as surviving (8/120).

The Metrics indicate:

- Accuracy: The number of records classified correctly by the Bagging Ensemble is 94 (66 + 27).
- %Accuracy: The percentage of records classified correctly is 78.33% (94/120).
- Specificity: (True Negative)/(True Negative + False Positives)  
 $68/(68 + 18) = 0.791$

Specificity is defined as the proportion of negative classifications that were actually negative, or the fraction of survivors that actually survived. In this model, 68 actual surviving patients were classified correctly as survivors. There were 18 false positives or 18 actual surviving patients were classified as deceased.

- Sensitivity or Recall: (True Positive)/(True Positive + False Negative)  
 $26/(26 + 8) = 0.765$

Sensitivity is defined as the proportion of positive cases there were classified correctly as positive, or the proportion of actually deceased patients there were classified as deceased. In this model, 26 actual deceased patients were correctly classified as deceased. There were 8

false negatives or 8 actual deceased patients were incorrectly classified as survivors.

- Precision: (True Positives)/(True Positives + False Positives)

$$26/(26 + 18) = 0.591$$

Precision is defined as the proportion of positive results that are true positive. In this model, 26 actual deceased patients were classified correctly as deceased. There were 18 false positives or 18 actual survivors classified incorrectly as deceased.

- F-1 Score:  $2 \times (\text{Precision} \times \text{Sensitivity}) / (\text{Precision} + \text{Sensitivity})$

$$2 \times (0.591 \times 0.765) / (0.591 + 0.765) = 0.667$$

The F-1 Score provides a statistic to balance between Precision and Sensitivity, especially if an uneven class distribution exists, as in this example, (86 survivors vs 34 deceased). The closer the F-1 score is to 1 (the upper bound) the better the precision and recall.

- Success Class and Success Probability simply reports the settings for these two values as input on the Find Best Model: Classification Data tab.

Individual records and their classifications are shown beneath Validation: Classification Details. Any misclassified records are highlighted in red.

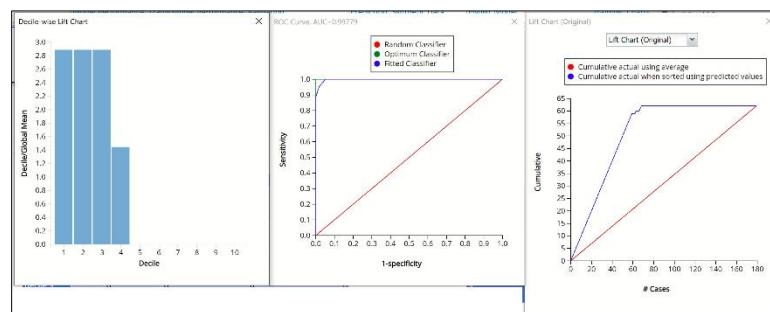
| Record ID  | DEATH_EVENT | Prediction: DEATH_EVENT | PostProb: 0 | PostProb: 1 |
|------------|-------------|-------------------------|-------------|-------------|
| Record 104 | 0           | 0                       | 0.6         | 0.4         |
| Record 163 | 0           | 0                       | 0.6         | 0.4         |
| Record 290 | 0           | 1                       | 0.5         | 0.5         |
| Record 207 | 0           | 0                       | 1           | 0           |
| Record 126 | 0           | 0                       | 0.9         | 0.1         |
| Record 228 | 0           | 0                       | 0.7         | 0.3         |

### CFBM\_TrainingLiftChart and CFBM\_ValidationLiftChart

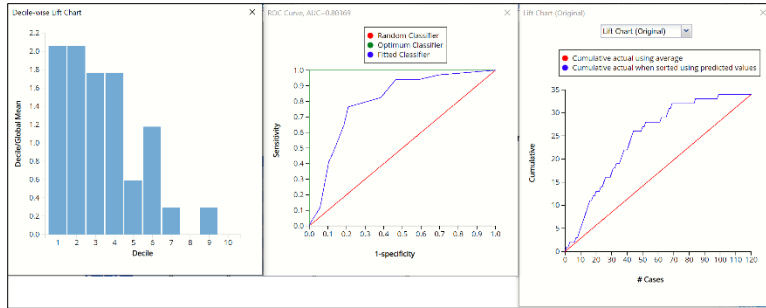
Lift Charts and ROC Curves are visual aids that help users evaluate the performance of their fitted models. Charts found on the CFBM\_TrainingLiftChart tab were calculated using the Training Data Partition. Charts found on the CFBM\_ValidationLiftChart tab were calculated using the Validation Data Partition. Since it is good practice to look at both sets of charts to assess model performance, this example will look at the lift charts generated from both partitions.

Note: To view these charts in the Cloud app, click the Charts icon on the Ribbon, select CFBM\_TrainingLiftChart or CFBM\_ValidationLiftChart.

### Decile-wise Lift Chart, ROC Curve, and Lift Charts for Training Partition



Decile-wise Lift Chart, ROC Curve, and Lift Charts for Valid. Partition



## Original Lift Chart

After the model is built using the training data set, the model is used to score on the training data set and the validation data set (if one exists). Afterwards, the data set(s) are sorted in decreasing order using the predicted output variable value. After sorting, the actual outcome values of the output variable are cumulated and the lift curve is drawn as the cumulative number of cases in decreasing probability (on the x-axis) vs the cumulative number of true positives on the y-axis. The baseline (red line connecting the origin to the end point of the blue line) is a reference line. For a given number of cases on the x-axis, this line represents the expected number of successes if no model existed, and instead cases were selected at random. This line can be used as a benchmark to measure the performance of the fitted model. The greater the area between the lift curve and the baseline, the better the model.

- In the Training Lift chart, if 100 records were selected as belonging to the success class and the fitted model was used to pick the members most likely to be successes, the lift curve indicates that 60 of them would be correct. Conversely, if 100 random records were selected and the fitted model was not used to identify the members most likely to be successes, the baseline indicates that about 30 would be correct.
- In the Validation Lift chart, if 100 records were selected as belonging to the success class and the fitted model was used to pick the members most likely to be successes, the lift curve indicates that about 40 classifications would be correct. Conversely, if 100 random records were selected and the fitted model was not used to identify the members most likely to be successes, the baseline indicates that about 25 would be correct.

## Decile-wise Lift Chart

The decilewise lift curve is drawn as the decile number versus the cumulative actual output variable value divided by the decile's mean output variable value. This bars in this chart indicate the factor by which the model outperforms a random assignment, one decile at a time.

- Refer to the decile-wise lift chart for the training dataset. In the first decile, the predictive performance of the model is almost 3 times better than simply assigning a random predicted value.
- Refer to the validation graph above. In the first decile, the predictive performance of the model is a little over 2 times better than simply assigning a random predicted value.

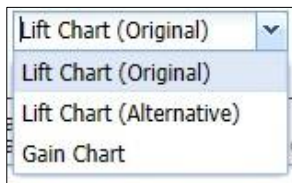
## ROC Curve

The Regression ROC curve was updated in V2017. This new chart compares the performance of the regressor (Fitted Predictor) with an Optimum Predictor Curve and a Random Classifier curve. The Optimum Predictor Curve plots a hypothetical model that would provide perfect classification results. The best possible classification performance is denoted by a point at the top left of the graph at the intersection of the x and y axis. This point is sometimes referred to as the “perfect classification”. The closer the AUC is to 1, the better the performance of the model.

- In the Training partition, AUC = .998 which suggests that this fitted model is a good fit to the data.
- In the Validation partition, AUC = .803 which suggests that this fitted model is also a good fit to the data.

## Alternative Lift Chart and Gain Chart

In V2017, two new charts were introduced: a new Lift Chart and the Gain Chart. To display these new charts, click the down arrow next to Lift Chart (Original), in the Original Lift Chart, then select the desired chart.

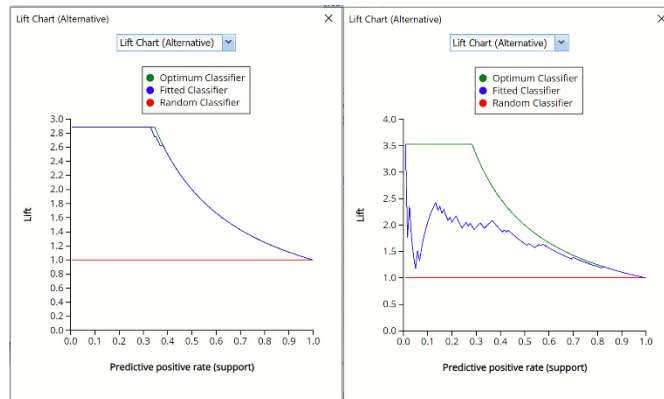


Select Lift Chart (Alternative) to display Analytic Solver Data Science's alternative Lift Chart. Each of these charts consists of an Optimum Predictor curve, a Fitted Predictor curve, and a Random Predictor curve.

The Optimum Predictor curve plots a hypothetical model that would provide perfect classification for our data. The Fitted Predictor curve plots the fitted model and the Random Predictor curve plots the results from using no model or by using a random guess (i.e. for x% of selected observations, x% of the total number of positive observations are expected to be correctly classified).

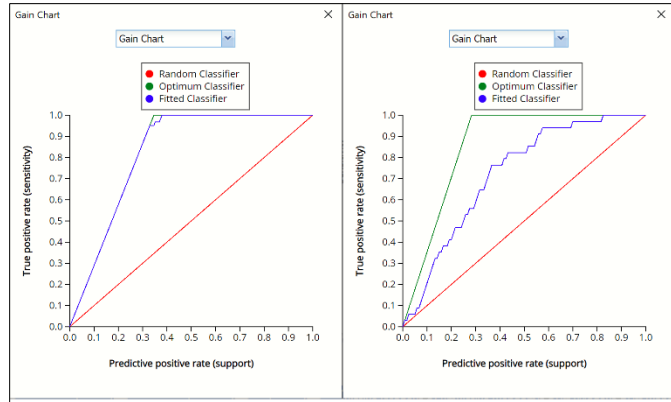
The Alternative Lift Chart plots Lift against the Predictive Positive Rate or Support.

### Lift Charts (Alternative) for Training and Validation Partitions



Click the down arrow and select Gain Chart from the menu. In this chart, the True Positive Rate or Sensitivity is plotted against the Predictive Positive Rate or Support.

### Gain Charts for Training and Validation Partitions



## CFBM\_Simulation

As discussed above, Analytic Solver Data Science generates a new output worksheet, CFBM\_Simulation, when *Simulate Response Prediction* is selected on the Simulation tab of the Find Best Model dialog.

This report contains the synthetic data, the predicted values for the training partition (using the fitted model) and the Excel – calculated Expression column, if populated, in the dialog. Users can switch between charts displaying the Predicted synthetic data, the Predicted training data or the Expression, or a combination of two, as long as they are of the same type.

### Synthetic Data

| Record   | Expression | DEATH_EVENT | age      | creatinine_phosphokinase | ejection_fraction | platelets | serum_creatinine | serum_sodium |
|----------|------------|-------------|----------|--------------------------|-------------------|-----------|------------------|--------------|
| Record 1 | EF>20      | 0           | 44.70336 | 515.0694613              | 42.9333681        | 72188.177 | 0.938832044      | 117.3417967  |
| Record 2 | EF>20      | 0           | 44.14701 | 1482.380512              | 46.35380581       | 127183.39 | 0.682460185      | 142.4089635  |
| Record 3 | EF>20      | 1           | 41.61868 | 461.9809875              | 36.94330059       | 179470.03 | 0.623244778      | 129.6897033  |
| Record 4 | EF>20      | 0           | 40.90336 | 852.2219761              | 58.11138604       | 344259.21 | 1.356385385      | 136.8318379  |
| Record 5 | EF>20      | 1           | 63.84847 | 1719.848596              | 26.61566          | 110998.94 | 1.084502545      | 120.1667257  |
| Record 6 | EF>20      | 1           | 40.00987 | 7711.39806               | 74.77711989       | 439671.61 | 0.611098249      | 138.2576534  |

Note the first column in the output, Expression. This column was inserted into the Synthetic Data results because Calculate Expression was selected and an Excel function was entered into the Expression field, on the Simulation tab of the Find Best Model dialog.

IF([@ejection\_fraction]<=20, [@DEATH\_EVENT], “EF > 20”)

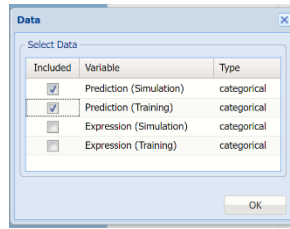
The results in this column are either 0, 1, or EF > 20.

- DEATH\_EVENT = 0 indicates that the patient has an ejection fraction <= 20 but *did not* suffer catastrophic heart failure.
- DEATH\_EVENT = 1 in this column indicates that the patient has an ejection fraction <= 20 and *did* suffer catastrophic heart failure.
- EF > 20 indicates that the patient’s ejection fraction was over 20.

The remainder of the data in this report is synthetic data, generated using the Generate Data feature described in the chapter with the same name, that appears earlier in this guide.

Click Prediction (Simulation), in the upper right hand corner of the chart, and select Prediction (Training) in the Data dialog to add the predicted data for the training partition into the chart.

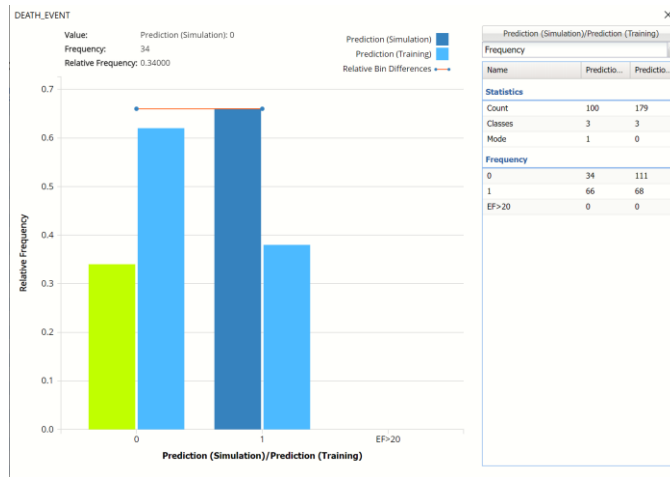
*Simulation Frequency Chart Data Dialog*



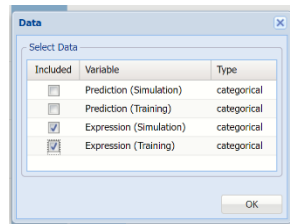
The chart that is displayed once this tab is selected, contains frequency information pertaining to the output variable’s predicted values in the synthetic data (darker shaded columns) and the training partition (lighter shaded columns) In the synthetic data, 34 patients are predicted to survive and 66 are not. In the training data, 111 patients were predicted to survive and 68 were not.

Note the flat Relative Bin Differences curve indicates the Absolute Difference in the bins is equal. (Click the down arrow next to Frequency and select Bin Details to find the bin metrics.)

*Frequency Chart for CFBM\_Simulation output*

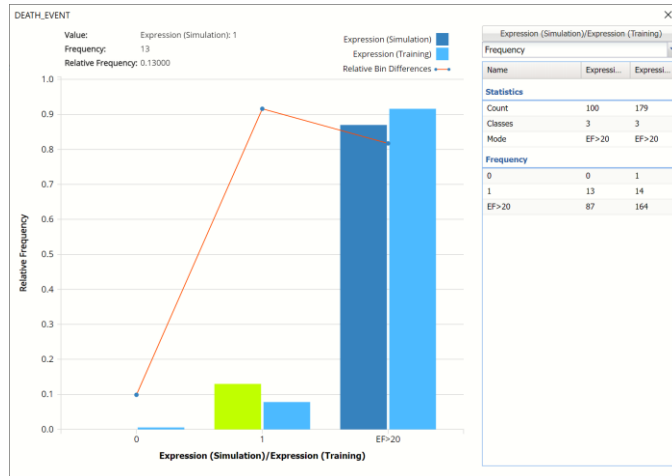


Click *Prediction(Simulation)/ Prediction (Training)* to change the chart view to *Expression (Simulation)/Expression (Training)* by selecting **Expression (Simulation)** and **Expression (Training)** in the Data dialog.



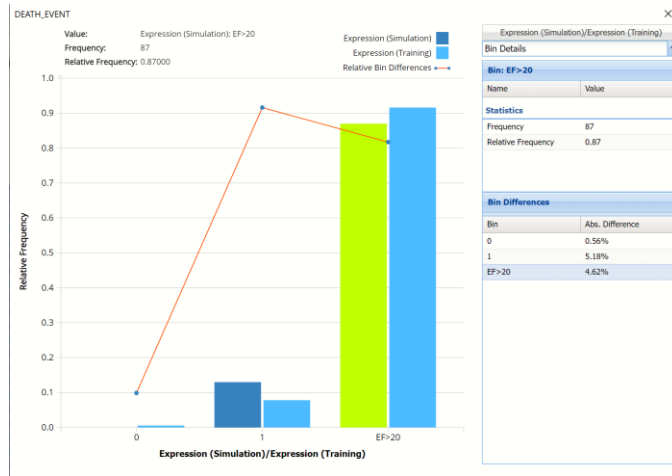
The chart below shows that in the synthetic, or simulated data, no “patients” with an ejection fraction  $\leq 20$  are expected to survive. However, in the training data, 1 patient with an ejection fraction  $\leq 20$  was predicted to survive. The 3<sup>rd</sup> column, EF>20 displays the remaining predictions for patients with ejection fraction  $> 20$  in both the simulation and training data.

### Frequency Chart for Expression



Click the down arrow next to Frequency and select Bin Details to see the Absolute Difference for each bin.

### Frequency Chart for Expression, Bin Details view



Notice that red lines which connect the relative Bin Differences for each bin. Bin Differences are computed based on the frequencies of records which predictions fall into each bin. For example, consider the highlighted bin in the screenshot above [Expression]. There are 87 Simulation records and 164 Training records in this bin. The relative frequency of the Simulation data is  $87/100 = 87\%$  and the relative frequency of the Training data is  $164/179 = 91.6\%$ . Hence the Absolute Difference (in frequencies) is  $91.6 - 87 = 4.6\%$

Click the down arrow next to Frequency to change the chart view to Relative Frequency or to change the look by clicking Chart Options. Statistics on the right of the chart dialog are discussed earlier in this section. For more information on the generated synthetic data, see the Generate Data chapter that appears earlier in this guide.

For information on Stored Model Sheets, in this example CFBM\_Stored, please refer to the “Scoring New Data” chapter within the Analytic Solver Data Science User Guide.



## Scoring New Data

Now that the model has been fit to the data, this fitted model will be used to score new patient data, found below.

### New Data

|   | A   | B       | C                        | D        | E                 | F                   | G         | H                | I            | J   | K       |
|---|-----|---------|--------------------------|----------|-------------------|---------------------|-----------|------------------|--------------|-----|---------|
| 6 | age | anaemia | creatinine_phosphokinase | diabetes | ejection_fraction | high_blood_pressure | platelets | serum_creatinine | serum_sodium | sex | smoking |
| 7 | 61  | 0       | 582                      | 0        | 38                | 0                   | 263358    | 1.4              | 137          | 1   | 0       |

To score new data, click the New Data worksheet tab and then click the Score icon on the Analytic Solver Data Science ribbon.

The screenshot shows the 'Scoring New Data' dialog box with several callouts:

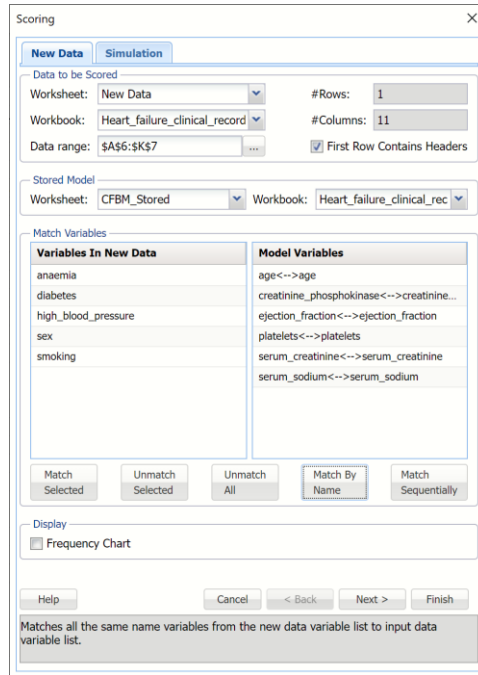
- Name of worksheet containing new data:** Points to the 'New Data' worksheet tab.
- Name of workbook:** Points to the 'Heart\_failure\_clinical\_record' workbook dropdown.
- Range for new data:** Points to the '\$A\$6:\$K\$7' data range field.
- Number of rows/columns in new data:** Points to the '# Rows: 1' and '# Columns: 11' fields.
- Variable names in new data:** Points to the 'Variables In New Data' list.
- Variable names in fitted model:** Points to the 'Model Variables' list.

The dialog box includes the following sections:

- Data to be Scored:** Worksheet (New Data), Workbook (Heart\_failure\_clinical\_record), Data range (\$A\$6:\$K\$7), # Rows (1), # Columns (11), and a checked 'First Row Contains Headers' option.
- Stored Model:** Worksheet (CFBM\_Stored), Workbook (Heart\_failure\_clinical\_rec).
- Match Variables:** Two lists: 'Variables In New Data' (age, anaemia, creatinine\_phosphokinase, diabetes, ejection\_fraction, high\_blood\_pressure, platelets, serum\_creatinine, serum\_sodium) and 'Model Variables' (age, creatinine\_phosphokinase, ejection\_fraction, platelets, serum\_creatinine, serum\_sodium). Buttons for 'Match Selected', 'Unmatch Selected', 'Unmatch All', 'Match By Name', and 'Match Sequentially' are present.
- Display:** A checked 'Frequency Chart' option.
- Buttons:** Help, Cancel, < Back, Next >, and Finish.

Click "Match By Name" to match each variable in the new data with the same variable in the fitted model, i.e. age with age, anaemia with anaemia, etc.

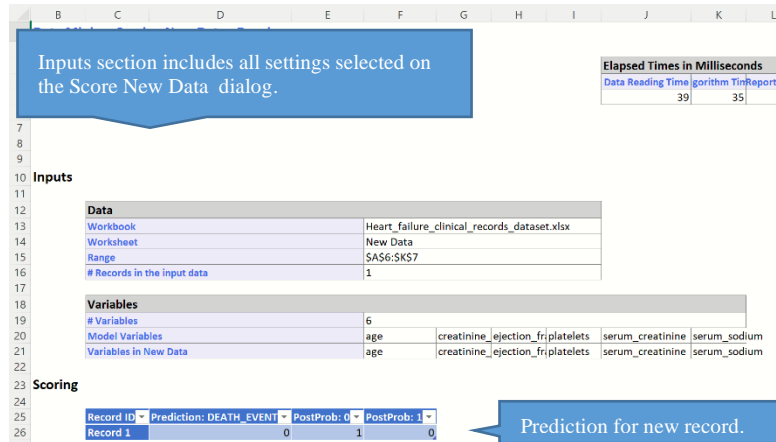
### Scoring tab



Click OK to score the new data record and determine if mitigation measures are needed to keep the patient healthy.

A new worksheet, Scoring\_Bagging is inserted to the right.

### Scoring Results



Notice that the predicted Death\_Event is 0, or that this patient will survive without any further mitigations at this point in time. Scoring should be repeated when changes in the patient's data is observed.

Please see the “Scoring New Data” chapter within the Analytic Solver Data Science User Guide for information on scoring new data.

## Find Best Model Regression Model

This example illustrates how to utilize the Find Best Model for Regression included in Analytic Solver Data Science for Desktop Excel or Excel Online by

using the Wine dataset<sup>7</sup>. This dataset contains 13 different features describing three wine varieties obtained from three different vineyards, all located in the same vicinity.

Find Best Model fits a model to all selected regression methods in order to observe which method provides the best fit to the data. The goal of this example is to fit the best model to the dataset, then use this fitted model to determine the alcohol content in a new sample of wine.

A list of each variable contained in the dataset appears in the table below.

| <b>VARIABLE</b>      |
|----------------------|
| Alcohol              |
| Malic_Acid           |
| Ash                  |
| Ash_Alcalinity       |
| Magnesium            |
| Total_Phenols        |
| Flavanoids           |
| Nonflavanoid_Phenols |
| Proanthocyanins      |
| Color_Intensity      |
| Hue                  |
| OD280_OD315          |
| Proline              |

All supervised algorithms include a new Simulation tab. This tab uses the functionality from the Generate Data feature (described in the What's New section of this guide and then more in depth in the Analytic Solver Data Science Reference Guide) to generate synthetic data based on the training partition, and uses the fitted model to produce predictions for the synthetic data. The resulting report, PFBM\_Simulation, will contain the synthetic data, the predicted values and the Excel-calculated Expression column, if present. In addition, frequency charts containing the Predicted, Training, and Expression (if present) sources or a combination of any pair may be viewed, if the charts are of the same type. Since this new functionality does not support categorical variables, these types of variables will not be present in the model, only continuous variables.

## Opening the Dataset

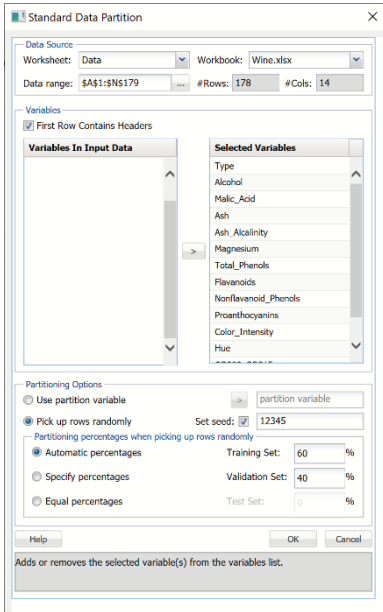
**Open wine.xlsx by clicking Help – Example Models – Forecasting/Data Science Examples.**

---

<sup>7</sup> This data set can be found in the UCI Machine Learning Repository (<http://www.ics.uci.edu/~mllearn/MLSummary.html> or <ftp://ftp.ics.uci.edu/pub/machine-learning-databases/wine/>)

## Partitioning the Dataset

Standard Data Partition Dialog



Partition the dataset by clicking **Partition – Standard Partition**.

- **Move all features under Variables In Input Data to Selected Variables.**
- Click **OK** to accept the partitioning defaults and create the partitions.

A new worksheet, *STDPartition*, is inserted directly to the right of the dataset. Click the new tab to open the worksheet. Find Best Model will be performed on both the Training and Validation partitions.

For more information on partitioning a dataset, see the Partitioning chapter within the Data Science Reference Guide.

## Running Find Best Model

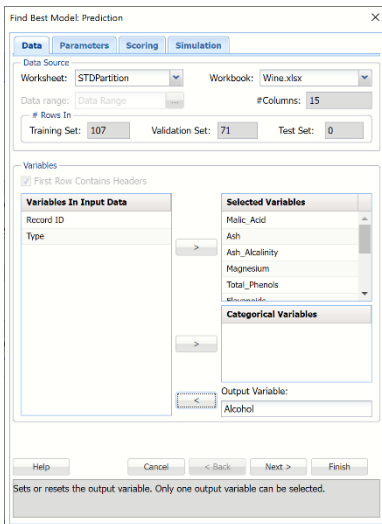
With *STDPartition* worksheet selected, click **Predict – Find Best Model** to open the Find Best Model Data tab.

### Data Tab

The continuous variables are selected on the Data tab.

- **Select Malic\_Acid, Ash, Ash\_Alcalinity, Magnesium, Total\_Phenols, Flavanoids, Nonflavanoid\_Phenols, Proanthocyanins, Color\_Intensity, Hue, OD280\_OD315, Proline for Selected Variables.**
- **Select Alcohol for the Output Variable.**
- Click **Next** to move to the Parameters tab.

Find Best Model: Prediction Data Tab



### Parameters tab

By default all eligible regression learners are automatically enabled based on the presence of categorical features or binary/multiclass classification. Optionally, all possible parameters for each algorithm may be defined using the Parameters button to the right of each learner.

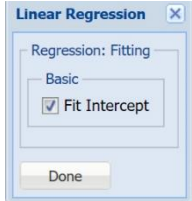
The Partition Data button is disabled because the original dataset was partitioned before Find Best Model was initiated.

### Rescale Data

To rescale the data, click Rescale Data and select the Rescale Data option at the top of the Rescaling dialog.

Use Rescaling to normalize one or more features in your data during the data preprocessing stage. Analytic Solver Data Science provides the following methods for feature scaling: Standardization, Normalization, Adjusted Normalization and Unit-norm.

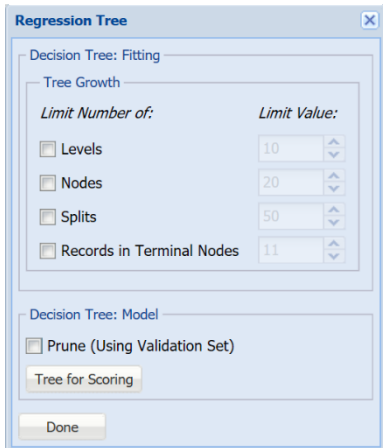
Linear Regression Parameter Dialog



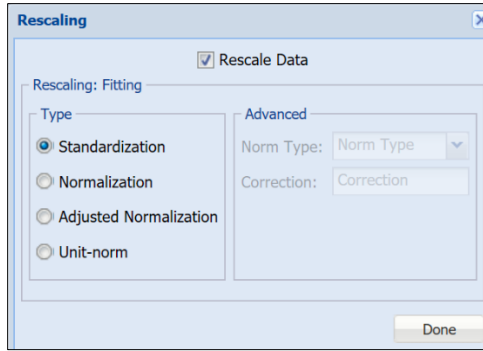
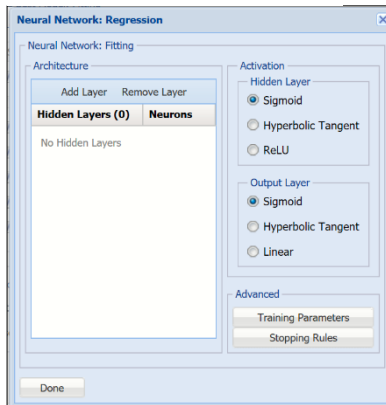
K-Nearest Neighbors Regression Parameter Dialog



Decision Trees Parameter Dialog



Neural Network Parameter Dialog



This example does not use rescaling to rescale the dataset data. Uncheck the Rescale Data option at the top of the dialog and click Done.

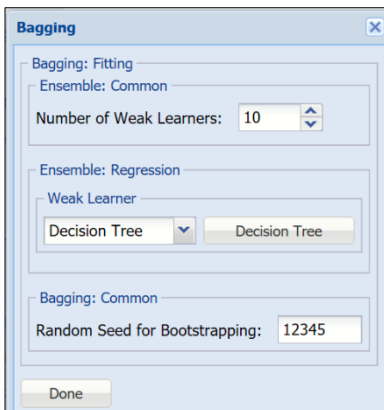
### Find Best Model: Fitting

To set a parameter for each selected learner, click the Parameters button to the right.

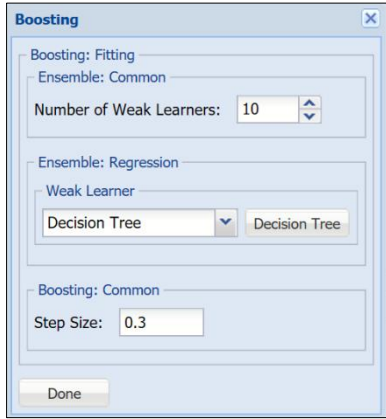
|                                                                   |                                                                                                                                                                                                                                                                                                                                                                                                                                             |
|-------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Linear Regression                                                 | For more information on each parameter, see the Linear Regression Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                   |
| Fit Intercept                                                     | If this option is selected, a constant term will be included in the model. Otherwise, a constant term will not be included in the equation. This option is selected by default.                                                                                                                                                                                                                                                             |
| K-Nearest Neighbors                                               | For more information on each parameter, see the K-Nearest Neighbors Regression Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                      |
| # Neighbors (k)                                                   | Enter a value for the parameter K in the Nearest Neighbor algorithm.                                                                                                                                                                                                                                                                                                                                                                        |
| Regression Tree                                                   | For more information on each parameter, see the Regression Tree Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                     |
| Tree Growth Levels, Nodes, Splits, Tree Records in Terminal Nodes | In the <i>Tree Growth</i> section, select Levels, Nodes, Splits, and Records in Terminal Nodes. Values entered for these options limit tree growth, i.e. if 10 is entered for Levels, the tree will be limited to 10 levels.                                                                                                                                                                                                                |
| Prune                                                             | If a validation partition exists, this option is enabled. When this option is selected, Analytic Solver Data Science will prune the tree using the validation set. Pruning the tree using the validation set reduces the error from over-fitting the tree to the training data.<br><br>Click Tree for Scoring to click the Tree type used for scoring: Fully Grown, Best Pruned, Minimum Error, User Specified or Number of Decision Nodes. |
| Neural Network                                                    | For more information on each parameter, see the Neural Network Regression Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                           |

|                               |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|-------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Architecture                  | Click <i>Add Layer</i> to add a hidden layer. To delete a layer, click <i>Remove Layer</i> . Once the layer is added, enter the desired Neurons.                                                                                                                                                                                                                                                                                                                                                                                   |
| Hidden Layer                  | Nodes in the hidden layer receive input from the input layer. The output of the hidden nodes is a weighted sum of the input values. This weighted sum is computed with weights that are initially set at random values. As the network “learns”, these weights are adjusted. This weighted sum is used to compute the hidden node’s output using a <i>transfer function</i> . The default selection is <i>Sigmoid</i> .                                                                                                            |
| Output Layer                  | As in the hidden layer output calculation (explained in the above paragraph), the output layer is also computed using the same transfer function as described for <i>Activation: Hidden Layer</i> . The default selection is <i>Sigmoid</i> .                                                                                                                                                                                                                                                                                      |
| Training Parameters           | Click Training Parameters to open the Training Parameters dialog to specify parameters related to the training of the Neural Network algorithm.                                                                                                                                                                                                                                                                                                                                                                                    |
| Stopping Rules                | Click Stopping Rules to open the Stopping Rules dialog. Here users can specify a comprehensive set of rules for stopping the algorithm early plus cross-validation on the training error.                                                                                                                                                                                                                                                                                                                                          |
| Bagging Ensemble Method       | For more information on each parameter, see the Ensemble Methods chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                  |
| Number of Weak Learners       | This option controls the number of “weak” regression models that will be created. The ensemble method will stop when the number of regression models created reaches the value set for this option. The algorithm will then compute the weighted sum of votes for each class and assign the “winning” value to each record.                                                                                                                                                                                                        |
| Weak Learner                  | Under Ensemble: Common click the down arrow beneath Weak Learner to select one of the four featured classifiers: Linear Regression, k-NN, Neural Networks or Decision Tree. The command button to the right will be enabled. Click this command button to control various option settings for the weak learner.                                                                                                                                                                                                                    |
| Random Seed for Bootstrapping | Enter an integer value to specify the seed for random resampling of the training data for each weak learner. Setting the random number seed to a nonzero value (any number of your choice is OK) ensures that the same sequence of random numbers is used each time the dataset is chosen for the classifier. The default value is “12345”. If left blank, the random number generator is initialized from the system clock, so the sequence of random numbers will be different in each calculation. If you need the results from |

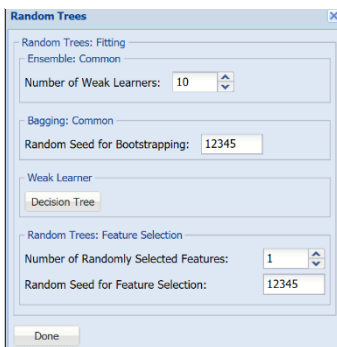
Ensemble Method Bagging Parameter Dialog



Ensemble Method Boosting Parameter Dialog



Ensemble Method Random Trees Parameter Dialog



|                                      |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
|--------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                      | successive runs of the algorithm to another to be strictly comparable, you should set the seed. To do this, type the desired number you want into the box. This option accepts positive integers with up to 9 digits.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| Boosting Ensemble Method             | For more information on each parameter, see the Boosting Regression Ensemble Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Number of Weak Learners              | See description above.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Weak Learner                         |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| Step Size                            | The Adaboost algorithm minimizes a loss function using the gradient descent method. The Step size option is used to ensure that the algorithm does not descend too far when moving to the next step. It is recommended to leave this option at the default of 0.3, but any number between 0 and 1 is acceptable. A Step size setting closer to 0 results in the algorithm taking smaller steps to the next point, while a setting closer to 1 results in the algorithm taking larger steps towards the next point.                                                                                                                                                                                                                                                                                                   |
| Random Trees Ensemble Method         | For more information on each parameter, see the Boosting Classification Ensemble Method chapter within the Data Science Reference Guide.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Number of Weak Learners              | See description above.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Random Seed for Bootstrapping        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| Weak Learner                         |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| Number of Randomly Selected Features | The Random Trees ensemble method works by training multiple “weak” classification trees using a fixed number of randomly selected features then taking the mode of each class to create a “strong” classifier. The option <i>Number of randomly selected features</i> controls the fixed number of randomly selected features in the algorithm. The default setting is <b>3</b> .                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Feature Selection Random Seed        | If an integer value appears for <i>Feature Selection Random seed</i> , Analytic Solver Data Science will use this value to set the feature selection random number seed. Setting the random number seed to a nonzero value (any number of your choice is OK) ensures that the same sequence of random numbers is used each time the dataset is chosen for the classifier. The default value is “12345”. If left blank, the random number generator is initialized from the system clock, so the sequence of random numbers will be different in each calculation. If you need the results from successive runs of the algorithm to another to be strictly comparable, you should set the seed. To do this, type the desired number you want into the box. This option accepts positive integers with up to 9 digits. |

**Find Best Model: Scoring**

The two parameters at the bottom of the dialog under Find Best Model: Scoring, determine how well each regression method fits the data.

The Metric for Scoring may be changed to R2, SSE, MSE, RMSE or MAD. See the table below for a brief description of each statistic.

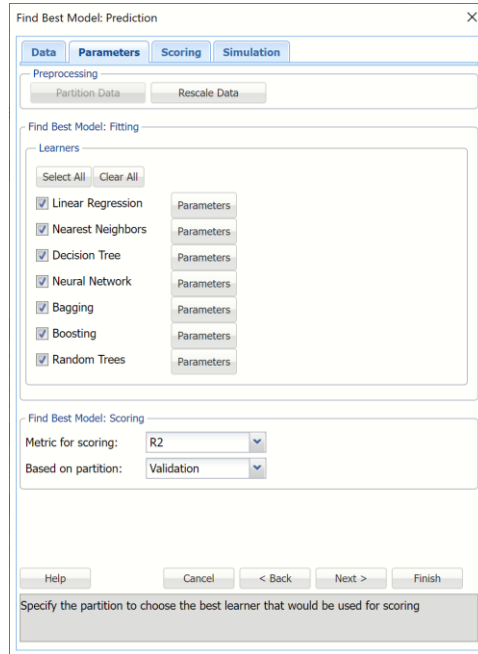
- For this example, **leave R2 selected**.
- Select **Validation** for Based on partition.

| Statistic | Description                                                                                                                                                                                                                                                                                                                                                                                                                                                          |
|-----------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| R2        | <p>Coefficient of Determination - Examines how differences in one variable can be explained by the difference in a second variable, when predicting the outcome of a given event.</p> $RMSE = \frac{SSR}{SST} = \frac{\sum_i(\hat{y}_i - \bar{y})^2}{\sum_i(y_i - \bar{y})^2}$ <p>where</p> <p><math>\hat{y}_i</math> is the predicted value for obs i</p> <p><math>y_i</math> is the actual value for obs i</p> <p><math>\bar{y}</math> is mean of the y values</p> |
| SSE       | <p>Sum of Squared Error – The sum of the squares of the differences between the actual and predicted values.</p> $SSE = \sum_{i=1}^n (\hat{y}_i - y_i)^2$ <p><math>\hat{y}_i</math> is the predicted value for obs i</p> <p><math>y_i</math> is the actual value for obs i</p>                                                                                                                                                                                       |
| MSE       | <p>Mean Squared Error – The average of the squared differences between the actual and predicted values.</p> $MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$ <p><math>\hat{y}_i</math> is the predicted value for obs i</p> <p><math>y_i</math> is the actual value for obs i</p>                                                                                                                                                                                |
| RMSE      | <p>Root Mean Squared Error – The standard deviation of the residuals.</p> $RSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$ <p><math>\hat{y}_i</math> is the predicted value for obs i</p> <p><math>y_i</math> is the actual value for obs i</p>                                                                                                                                                                                                             |
| MAD       | <p>Mean Absolute Deviation - Average distance between each data value and the sample mean; describes variation in a data set.</p>                                                                                                                                                                                                                                                                                                                                    |



|  |                                                                                                                                                                                       |
|--|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|  | $MAD = \frac{1}{n} \sum_{i=1}^n  x_i - \bar{x} $ <p>where <math>x_i</math> is the <math>i^{\text{th}}</math> obs in the sample<br/> where <math>\bar{x}</math> is the sample mean</p> |
|--|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

Find Best Model dialog, Parameters tab



### Scoring tab

Output options are selected on the Scoring tab. By default, the PFBM\_Output worksheet will be inserted directly to the right of the STDPartition worksheet.

In this example select all four options: Detailed Report, Summary Report, Lift Charts and Frequency Chart under both Score Training Data and Score Validation Data. Then click Finish to run Find Best Model.

- PFBM\_Output contains a listing of all model inputs such as input/output variables and parameter settings for all Learners, as well as Model Performance tables containing evaluations for every available metric, every learner on all available partitions. The Learner identified internally as performing the best, is highlighted in red. (Recall that the statistic used for determining which Learner performs best on the dataset was selected on the Parameters tab.)
- PFBM\_Stored contains the PMML (Predictive Model Markup Language) model which can be utilized to score new data. For more information on scoring, see the Scoring chapter that appears later in this guide.

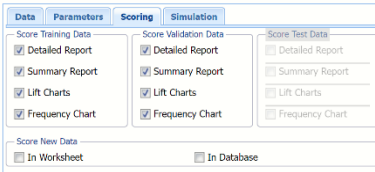
Selecting *Detailed Report* produces PFBM\_TrainingScore and PFBM\_ValidationScore.

- Both reports contain detailed scoring information on both the Training and Validation partitions using the "best" learner.

*Summary Report* is selected by default. This option produces a summary report at the top of both PFBM\_TrainingScore and PFBM\_ValidationScore worksheets.

- *Summary Report* contains a listing of 5 metrics: SSE, MSE, RMSE, MAD and R2.

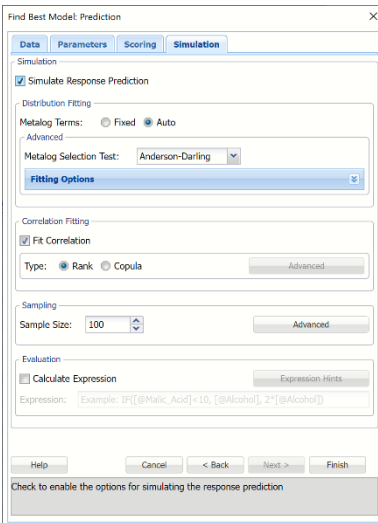
Find Best Model dialog, Scoring tab



Selecting *Frequency Chart* produces a frequency graph of the records in the partitions.

- When *Frequency Chart* is selected, a frequency chart will be displayed when the PFBM\_TrainingScore and PFBM\_ValidationScore worksheets are selected. This chart will display an interactive application similar to the Analyze Data feature, explained in detail in the Analyze Data chapter that appears earlier in this guide. This chart will include frequency distributions of the actual and predicted responses individually, or side-by-side, depending on the user's preference, as well as basic and advanced statistics for variables, percentiles, six sigma indices.

Find Best Model dialog, Simulation tab



Selecting *Lift Charts* generates Lift Charts, RROC Curves and Decil-Wise Lift Charts.

See the Scoring chapter that appears later in this guide for more information on the Score New Data section of the Scoring tab.

## Simulation Tab

Click **Next** to advance to the Simulation tab.

Select Simulation Response Prediction to enable all options on the Simulation tab of the Find Best Model Prediction dialog.

**Simulation tab:** All supervised algorithms include a new Simulation tab. This tab uses the functionality from the Generate Data feature (described earlier in this guide) to generate synthetic data based on the training partition, and uses the fitted model to produce predictions for the synthetic data. The resulting report, PFBM\_Simulation, will contain the synthetic data, the predicted values and the Excel-calculated Expression column, if present. In addition, frequency charts containing the Predicted, Training, and Expression (if present) sources or a combination of any pair may be viewed, if the charts are of the same type.

**Evaluation:** If Calculate Expression is selected, Analytic Solver amends an Expression column onto the frequency chart displayed on the PFBM\_Simulation output tab. Expression can be any valid Excel formula that references a variable and the response as [@COLUMN\_NAME]. Click the *Expression Hints* button for more information on entering an expression.

For the purposes of this example, leave this checkbox *unchecked*. See the Find Best Model classification example above to see an example of this option in use.

For more information on the remaining options shown on this dialog in the Distribution Fitting, Correlation Fitting and Sampling sections, see the Generate Data chapter that appears earlier in this guide.

Click **Finish** to run Find Best Model on the example dataset.

## Interpreting the Results

All output sheets are inserted to the right of STDPartition.

## PFBM\_Output

The PFBM\_Output worksheet is inserted directly to the right of the STDPartition worksheet. This report lists all input variables and all parameter settings for each learner, along with the Model Performance of each Learner on all partitions. This example utilizes two partitions, training and validation.

The Output Navigator appears at the very top of all output worksheets. Click the links to easily move between each section of the output. The Output Navigator is listed at the top of each worksheet included in the output.

|   |                                              |                                              |                                               |                                                |                                                |                                  |  |  |  |  |  |
|---|----------------------------------------------|----------------------------------------------|-----------------------------------------------|------------------------------------------------|------------------------------------------------|----------------------------------|--|--|--|--|--|
| 1 | Data Science: Find Best Model Prediction     |                                              |                                               |                                                |                                                |                                  |  |  |  |  |  |
| 2 |                                              |                                              |                                               |                                                |                                                |                                  |  |  |  |  |  |
| 3 | Output Navigator                             |                                              |                                               |                                                |                                                |                                  |  |  |  |  |  |
| 4 | <a href="#">Inputs</a>                       | <a href="#">Model Performance: Training</a>  | <a href="#">Model Performance: Validation</a> | <a href="#">Prediction: Synthetic Data</a>     | <a href="#">PMMI Model</a>                     | <a href="#">Training: Charts</a> |  |  |  |  |  |
| 5 | <a href="#">Training: Prediction Summary</a> | <a href="#">Training: Prediction Details</a> | <a href="#">Validation: Charts</a>            | <a href="#">Validation: Prediction Summary</a> | <a href="#">Validation: Prediction Details</a> |                                  |  |  |  |  |  |

The Inputs section includes information pertaining to the dataset, the input variables and parameter settings.

|    |                                              |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
|----|----------------------------------------------|------------|-----|----------------|-----------|---------------|------------|----------------------|-----------------|-----------------|-------------------|------------------------------------------------|---------|
| 12 | Data                                         |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 13 | Workbook                                     |            |     |                |           |               |            |                      |                 |                 | Wine.xlsx         |                                                |         |
| 14 | Worksheet                                    |            |     |                |           |               |            |                      |                 |                 | STDPartition      |                                                |         |
| 15 | Training data used for building the model    |            |     |                |           |               |            |                      |                 |                 | \$C\$37:\$Q\$143  |                                                |         |
| 16 | # Records in the training data               |            |     |                |           |               |            |                      |                 |                 | 107               |                                                |         |
| 17 | Validation data                              |            |     |                |           |               |            |                      |                 |                 | \$C\$144:\$Q\$214 |                                                |         |
| 18 | # Records in the validation data             |            |     |                |           |               |            |                      |                 |                 | 71                |                                                |         |
| 19 |                                              |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 20 | Variables                                    |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 21 | # Variables                                  |            |     |                |           |               |            |                      |                 |                 |                   | 12                                             |         |
| 22 | Scale Variables                              | Malic_Acid | Ash | Ash_Alcalinity | Magnesium | Total_Phenols | Flavanoids | Nonflavanoid_Phenols | Proanthocyanins | Color_Intensity | Hue               | OD280_OD31                                     | Proline |
| 23 | Categorical Variables                        |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 24 | Output Variable                              |            |     |                |           |               |            |                      |                 |                 |                   | Alcohol                                        |         |
| 25 |                                              |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 26 | Rescaling: Fitting Parameters                |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 27 | Rescale Data?                                |            |     |                |           |               |            |                      |                 |                 |                   | TRUE                                           |         |
| 28 | Technique                                    |            |     |                |           |               |            |                      |                 |                 |                   | STANDARDIZATION                                |         |
| 29 |                                              |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 30 | Find Best Model: Fitting Parameters          |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 31 | Learners                                     |            |     |                |           |               |            |                      |                 |                 |                   | LINEAR_REGRESSION, NEAREST_NEIGHBORS, DECISION |         |
| 32 |                                              |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 33 | Find Best Model Prediction: Model Parameters |            |     |                |           |               |            |                      |                 |                 |                   |                                                |         |
| 34 | Metric for scoring                           |            |     |                |           |               |            |                      |                 |                 |                   | R2                                             |         |
| 35 | Based on partition                           |            |     |                |           |               |            |                      |                 |                 |                   | Validation                                     |         |

Further down within Inputs, the parameter selections for each Learner are listed.

37 **Linear Regression**

| Regression Model: Fitting Parameters |      |
|--------------------------------------|------|
| Fit intercept                        | TRUE |

42 **Nearest Neighbors**

| Nearest Neighbors: Fitting Parameters |   |
|---------------------------------------|---|
| # Nearest neighbors (K)               | 1 |

47 **Decision Tree**

| Decision Tree: Model Parameters |             |
|---------------------------------|-------------|
| Prune?                          | FALSE       |
| Scoring tree type               | Fully grown |

53 **Neural Network**

| Neural Network: Fitting Parameters |                  |
|------------------------------------|------------------|
| Random seed for initial weights    | 12345            |
| # Hidden Layers                    | 0                |
| Learning rate                      | 0.1              |
| Weight change momentum             | 0.6              |
| Error tolerance                    | 0.01             |
| Weight decay                       | 0                |
| Cost function                      | Sum of Squares   |
| Hidden layer activation function   | LOGISTIC SIGMOID |
| Output layer activation function   | LOGISTIC SIGMOID |
| Learning order                     | Original         |
| Response correction                | 0.01             |
| Data for error computation         | TRAINING ONLY    |
| Maximum number of epochs           | 30               |
| Maximum number of epochs with      | 5                |
| Maximum training time              | 3600             |
| Minimum relative change in error   | 0.0001           |
| Minimum relative change in error   | 0.001            |

74 **Bagging**

| Ensemble Parameters       |                 |
|---------------------------|-----------------|
| Weak learner              | Regression Tree |
| Number of weak learners   | 10              |
| Show weak learner models? | FALSE           |

| Bagging Parameters |       |
|--------------------|-------|
| Bootstrap seed     | 12345 |

84 **Boosting**

| Ensemble Parameters       |                 |
|---------------------------|-----------------|
| Weak learner              | Regression Tree |
| Number of weak learners   | 10              |
| Show weak learner models? | FALSE           |

| Boosting Regression: Fitting Parameters |     |
|-----------------------------------------|-----|
| Step size                               | 0.3 |

94 **Random Trees**

| Ensemble Parameters       |                 |
|---------------------------|-----------------|
| Weak learner              | Regression Tree |
| Number of weak learners   | 10              |
| Show weak learner models? | FALSE           |

| Bagging Parameters |       |
|--------------------|-------|
| Bootstrap seed     | 12345 |

| Random Trees: Fitting Parameters |       |
|----------------------------------|-------|
| # Selected features              | 3     |
| Feature selection random seed    | 12345 |

Scroll down to view Simulation tab option settings and any generated messages from the Find Best Model feature.

Further down, the Model Performance tables display how each prediction method performed on the dataset.

|     | B | C                                                  | D                                               | E | F | G | H | I |
|-----|---|----------------------------------------------------|-------------------------------------------------|---|---|---|---|---|
| 107 |   | <b>Simulation: Distribution Fitting Parameters</b> |                                                 |   |   |   |   |   |
| 108 |   | MetaLog Terms                                      | Auto                                            |   |   |   |   |   |
| 109 |   | GOF Test                                           | Anderson-Darling                                |   |   |   |   |   |
| 110 |   | Options                                            | {"Malic_Acid":{"numTerms":5},"Ash":{"numTerms": |   |   |   |   |   |
| 111 |   | <b>Simulation: Correlation Fitting Parameters</b>  |                                                 |   |   |   |   |   |
| 112 |   | Correlation Type                                   | Rank                                            |   |   |   |   |   |
| 113 |   | <b>Simulation: Sampling Parameters</b>             |                                                 |   |   |   |   |   |
| 114 |   | Generate sample                                    | Yes                                             |   |   |   |   |   |
| 115 |   | Sample size                                        | 100                                             |   |   |   |   |   |
| 116 |   | Random seed                                        | 12345                                           |   |   |   |   |   |
| 117 |   | Random generator                                   | Mersenne Twister                                |   |   |   |   |   |
| 118 |   | Sampling method                                    | Latin Hypercube                                 |   |   |   |   |   |
| 119 |   | Random streams                                     | Independent                                     |   |   |   |   |   |
| 120 |   | Calculate expression?                              | No                                              |   |   |   |   |   |
| 121 |   | <b>Output Options</b>                              |                                                 |   |   |   |   |   |
| 122 |   | Summary report of scoring on training data         |                                                 |   |   |   |   |   |
| 123 |   | Detailed report of scoring on training data        |                                                 |   |   |   |   |   |
| 124 |   | Lift charts on training data                       |                                                 |   |   |   |   |   |
| 125 |   | Frequency chart on training data                   |                                                 |   |   |   |   |   |
| 126 |   | Summary report of scoring on validation data       |                                                 |   |   |   |   |   |
| 127 |   | Detailed report of scoring on validation data      |                                                 |   |   |   |   |   |
| 128 |   | Lift charts on validation data                     |                                                 |   |   |   |   |   |
| 129 |   | Frequency chart on validation data                 |                                                 |   |   |   |   |   |
| 130 |   | <b>Note:</b> Scoring will be done using Bagging    |                                                 |   |   |   |   |   |

The Messages portion of the report indicates that Scoring will be performed using the Bagging Ensemble Method the Learner selected as the "best" choice according to the selection for Find Best Model: Scoring parameters on the Parameters tab: Validation Partition R2 Metric.

#### 138 Model Performance: Training

| Metric            | SSE      | MSE         | RM   | MAD         | R2        |
|-------------------|----------|-------------|------|-------------|-----------|
| Linear Regression | 23.80754 | 0.222500349 | 0.47 | 0.383264031 | 0.639702  |
| Decision Tree     | 0        | 0           | 0    | 0           | 1         |
| Nearest Neighbors | 0        | 0           | 0    | 0           | 1         |
| Neural Network    | 75.76374 | 0.708072297 | 0.84 | 0.650156797 | -0.146593 |
| Bagging           | 4.408646 | 0.041202299 | 0.2  | 0.156691589 | 0.93328   |
| Boosting          | 0.052725 | 0.000492753 | 0.02 | 0.019134102 | 0.999202  |
| Random Trees      | 4.408646 | 0.041202299 | 0.2  | 0.156691589 | 0.93328   |

#### 149 Model Performance: Validation

| Metric            | SSE      | MSE         | RM   | MAD         | R2       |
|-------------------|----------|-------------|------|-------------|----------|
| Linear Regression | 26.62183 | 0.374955371 | 0.61 | 0.475844018 | 0.443059 |
| Decision Tree     | 43.7408  | 0.616067606 | 0.78 | 0.643661972 | 0.084922 |
| Nearest Neighbors | 33.8716  | 0.477064789 | 0.69 | 0.52028169  | 0.291391 |
| Neural Network    | 69.91289 | 0.984688548 | 0.99 | 0.794157969 | -0.46261 |
| Bagging           | 26.25256 | 0.36975431  | 0.61 | 0.480788732 | 0.450784 |
| Boosting          | 42.56739 | 0.599540748 | 0.77 | 0.634686457 | 0.109471 |
| Random Trees      | 26.25256 | 0.36975431  | 0.61 | 0.480788732 | 0.450784 |

Since the Bagging Ensemble R2 metric for the Validation Partition has the highest score, that is the Learner that will be used for scoring.

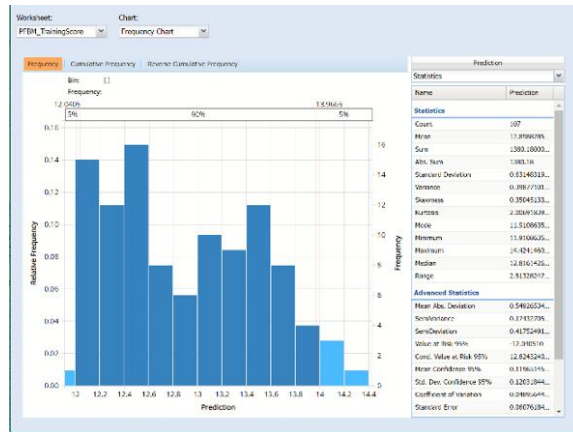
### ***PFBM\_TrainingScore and PFBM\_ValidationScore***

PFBM\_TrainingScore contains the Prediction Summary and the Prediction Details reports for the training partition. PFBM\_Validation contains the same reports for the validation partition. Both reports have been generated using the Bagging Ensemble Method, as discussed above.

#### **PFBM\_TrainingScore**

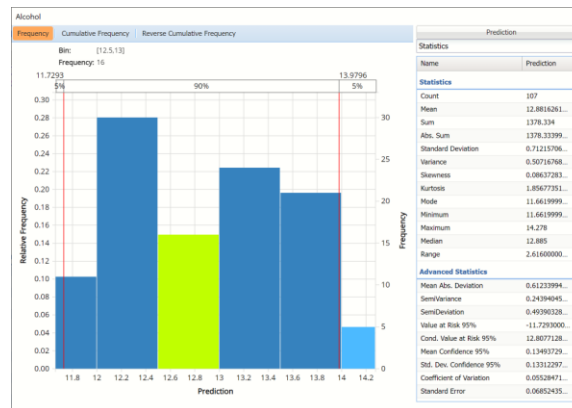
Click the *PFBM\_TrainingScore* tab to view the newly added Output Variable frequency chart, the Training: Prediction Summary and the Training: Prediction Details report. All calculations, charts and predictions on this worksheet apply to the Training data.

Note: To view charts in the Cloud app, click the Charts icon on the Ribbon, select a worksheet under Worksheet and a chart under Chart.



**Tabs:** The Analyze Data dialog contains three tabs: Frequency, Cumulative Frequency, and Reverse Cumulative Frequency. Each tab displays different information about the distribution of variable values.

### Analyze Data Dialog



Hovering over a bar in either of the three charts will populate the Bin and Frequency headings at the top of the chart. In the Frequency chart above, the bar for the [12.5, 13] Bin is selected. This bar has a frequency of 16 and a relative frequency of about 15%.

By default, red vertical lines will appear at the 5% and 95% percentile values in all three charts, effectively displaying the 90<sup>th</sup> confidence interval. The middle percentage is the percentage of all the variable values that lie within the ‘included’ area, i.e. the darker shaded area. The two percentages on each end are the percentage of all variable values that lie outside of the ‘included’ area or the “tails”. i.e. the lighter shaded area. Percentile values can be altered by moving either red vertical line to the left or right.

Click the “X” in the upper right corner of the detailed chart dialog to close the dialog. To re-open the chart, click a new tab, say the Data tab in this example, and then click PFBM\_TrainingScore.

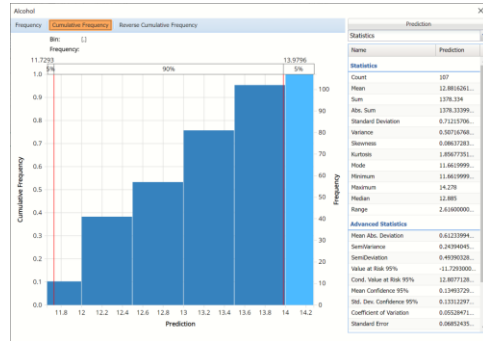
**Frequency Tab:** When the Analyze Data dialog is first displayed, the Frequency tab is selected by default. This tab displays a histogram of the variable’s values.

Bins containing the range of values for the variable appear on the horizontal axis, the relative frequency of occurrence of the bin values appears on the left vertical axis while the actual frequency of the bin values appear on the right vertical axis.

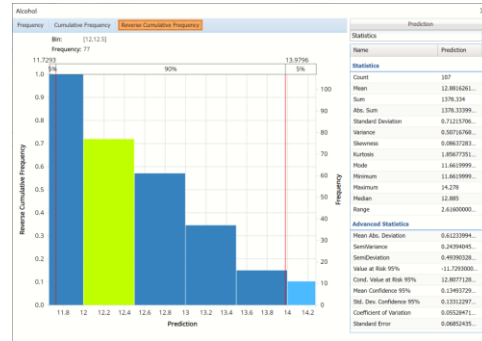
## Cumulative Frequency / Reverse Cumulative Frequency

The Cumulative Frequency tab displays a chart of the cumulative form of the frequency chart, as shown below. Hover over each bar to populate the Bin and Frequency headings at the top of the chart. In this screenshot below, the bar for the [12.5, 13.0] Bin is selected in the Cumulative Frequency Chart. This bar has a frequency of 57 and a relative frequency of about 52%.

Cumulative Frequency Chart



Reverse Cumulative Frequency Chart

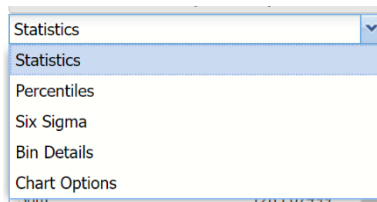


**Cumulative Frequency Chart:** Bins containing the range of values for the variable appear on the horizontal axis, the cumulative frequency of occurrence of the bin values appear on the left vertical axis while the actual cumulative frequency of the bin values appear on the right vertical axis.

**Reverse Cumulative Frequency Chart:** Bins containing the range of values for the variable appear on the horizontal axis, similar to the Cumulative Frequency Chart. The reverse cumulative frequency of occurrence of the bin values appear on the left vertical axis while the actual reverse cumulative frequency of the bin values appear on the right vertical axis.

Click the drop down menu on the upper right of the dialog to display additional panes: Statistics, Six Sigma and Percentiles.

Drop down menu



### Statistics View

The Statistics tab displays numeric values for several summary statistics, computed from all values for the specified variable. The statistics shown on the pane below were computed for the predicted variable, in this case Alcohol.

### Statistics Pane

| Name                       | Prediction     |
|----------------------------|----------------|
| <b>Statistics</b>          |                |
| Count                      | 107            |
| Mean                       | 12.8816261...  |
| Sum                        | 1378.334       |
| Abs. Sum                   | 1378.33399...  |
| Standard Deviation         | 0.71215706...  |
| Variance                   | 0.50716768...  |
| Skewness                   | 0.08637283...  |
| Kurtosis                   | 1.85677351...  |
| Mode                       | 11.6619999...  |
| Minimum                    | 11.6619999...  |
| Maximum                    | 14.278         |
| Median                     | 12.885         |
| Range                      | 2.6160000...   |
| <b>Advanced Statistics</b> |                |
| Mean Abs. Deviation        | 0.61233994...  |
| SemiVariance               | 0.24394045...  |
| SemiDeviation              | 0.49390328...  |
| Value at Risk 95%          | -11.7293000... |
| Cond. Value at Risk 95%    | 12.8077128...  |
| Mean Confidence 95%        | 0.13493729...  |
| Std. Dev. Confidence 95%   | 0.13312297...  |
| Coefficient of Variation   | 0.05528471...  |
| Standard Error             | 0.06852435...  |

All statistics appearing on the Statistics pane are briefly described below.

#### Statistics

- **Count**, the total number of records included in the data.
- **Mean**, the average of all the values.
- **Sum**, the sum of all the values.
- **Abs. Sum**, the absolute sum of all the values.
- **Standard Deviation**, the square root of variance.
- **Variance**, describes the spread of the distribution of values.
- **Skewness**, which describes the *asymmetry* of the distribution of values.
- **Kurtosis**, which describes the *peakedness* of the distribution of values.
- **Mode**, the most frequently occurring single value.
- **Minimum**, the minimum value attained.
- **Maximum**, the maximum value attained.
- **Median**, the median value.
- **Range**, the difference between the maximum and minimum values.

#### Advanced Statistics

- **Mean Abs. Deviation**, returns the average of the absolute deviations.
- **SemiVariance**, measure of the dispersion of values.



- **SemiDeviation**, *one-sided* measure of dispersion of values.
- **Value at Risk 95%**, the maximum loss that can occur at a given confidence level.
- **Cond. Value at Risk**, is defined as the *expected value* of a loss *given that* a loss at the specified percentile occurs.
- **Mean Confidence 95%**, returns the confidence “half-interval” for the estimated mean value (returned by the PsiMean() function).
- **Std. Dev. Confidence 95%**, returns the confidence ‘half-interval’ for the estimated standard deviation of the simulation trials (returned by the PsiStdDev() function).
- **Coefficient of Variation**, is defined as the ratio of the standard deviation to the mean.
- **Standard Error**, defined as the standard deviation of the sample mean.
- **Expected Loss**, returns the average of all negative data multiplied by the percentrank of 0 among all data.
- **Expected Loss Ratio**, returns the expected loss ratio.
- **Expected Gain** returns the average of all positive data multiplied by 1 - percentrank of 0 among all data.
- **Expected Gain Ratio**, returns the expected gain ratio.
- **Expected Value Margin**, returns the expected value margin.

### Percentiles View

Selecting Percentiles from the menu displays numeric percentile values (from 1% to 99%) computed using all values for the variable. The percentiles shown below were computed using the values for the Alcohol variable.

Percentiles Pane

| Prediction  |               |
|-------------|---------------|
| Percentiles |               |
| Name        | Prediction    |
| 82%         | 13.7011200... |
| 83%         | 13.7265       |
| 84%         | 13.7392       |
| 85%         | 13.7447       |
| 86%         | 13.751        |
| 87%         | 13.7512200... |
| 88%         | 13.75592      |
| 89%         | 13.78402      |
| 90%         | 13.8222       |
| 91%         | 13.8550600... |
| 92%         | 13.9004800... |
| 93%         | 13.94042      |
| 94%         | 13.9648399... |
| 95%         | 13.9795999... |
| 96%         | 14.0625199... |
| 97%         | 14.09438      |
| 98%         | 14.11712      |
| 99%         | 14.1848599... |

The values displayed here represent 99 equally spaced points on the Cumulative

Frequency chart: In the Percentile column, the numbers rise smoothly on the vertical axis, from 0 to 1.0, and in the Value column, the corresponding values from the horizontal axis are shown. For example, the 75th Percentile value is a number such that three-quarters of the values occurring in the last simulation are less than or equal to this value.

## **Six Sigma View**

Selecting Six Sigma from the menu displays various computed Six Sigma measures. In this display, the red vertical lines on the chart are the Lower Specification Limit (LSL) and the Upper Specification Limit (USL) which are initially set equal to the 5<sup>th</sup> and 95<sup>th</sup> percentile values, respectively.

These functions compute values related to the Six Sigma indices used in manufacturing and process control. For more information on these functions, see the Appendix located at the end of the Data Science Reference Guide.

- **SigmaCP** calculates the Process Capability.
- **SigmaCPK** calculates the Process Capability Index.
- **SigmaCPKLower** calculates the one-sided Process Capability Index based on the Lower Specification Limit.
- **SigmaCPKUpper** calculates the one-sided Process Capability Index based on the Upper Specification Limit.
- **SigmaCPM** calculates the Taguchi Capability Index.
- **SigmaDefectPPM** calculates the Defect Parts per Million statistic.
- **SigmaDefectShiftPPM** calculates the Defective Parts per Million statistic with a Shift.
- **SigmaDefectShiftPPMLower** calculates the Defective Parts per Million statistic with a Shift below the Lower Specification Limit.
- **SigmaDefectShiftPPMUpper** calculates the Defective Parts per Million statistic with a Shift above the Upper Specification Limit.
- **SigmaK** calculates the Measure of Process Center.
- **SigmaLowerBound** calculates the Lower Bound as a specific number of standard deviations below the mean.
- **SigmaProbDefectShift** calculates the Probability of Defect with a Shift outside the limits.
- **SigmaProbDefectShiftLower** calculates the Probability of Defect with a Shift below the lower limit.
- **SigmaProbDefectShiftUpper** calculates the Probability of Defect with a Shift above the upper limit.
- **SigmaSigmaLevel** calculates the Process Sigma Level with a Shift.
- **SigmaUpperBound** calculates the Upper Bound as a specific number of standard deviations above the mean.
- **SigmaYield** calculates the Six Sigma Yield with a shift, i.e. the fraction of the process that is free of defects.
- **SigmaZLower** calculates the number of standard deviations of the process that the lower limit is below the mean of the process.

- **SigmaZMin** calculates the minimum of ZLower and ZUpper.
- **SigmaZUpper** calculates the number of standard deviations of the process that the upper limit is above the mean of the process.

*Six Sigma Pane*

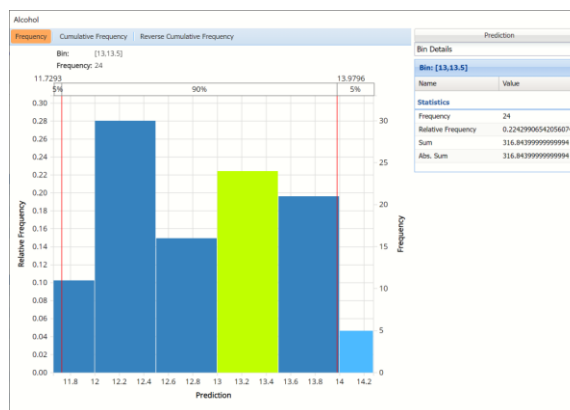
| Prediction      |                |
|-----------------|----------------|
| Six Sigma       |                |
| Name            | Prediction     |
| Cpk             | 0.51405452...  |
| Cpk, lower      | 0.53937592...  |
| Cpk, upper      | 0.51405452...  |
| Cpm             | 0.02907102...  |
| PPM             | 114334.458...  |
| PPM, lower      | 52817.53314... |
| PPM, upper      | 61516.9251...  |
| K               | 0.02403708...  |
| Lower Bound     | 8.60916769...  |
| ProbDefectShift | 0.11433445...  |
| ProbLowerShift  | 0.05281753...  |
| ProbUpperShift  | 0.06151692...  |
| Sigma Level     | 1.20379475...  |
| Upper Bound     | 17.15382296... |
| Yield           | 0.88566554...  |
| Z Lower         | 1.61812776...  |
| Z Min           | 1.54216356...  |
| Z Upper         | 1.54216356...  |

**Bin Details View**

Click the down arrow next to Statistics to view Bin Details for each bin in the chart.

**Bin:** If viewing the chart with only the Predicted or simulate data, only one grid will be displayed on the Bin Details pane. This grid displays important bin statistics such as frequency, relative frequency, sum and absolute sum.

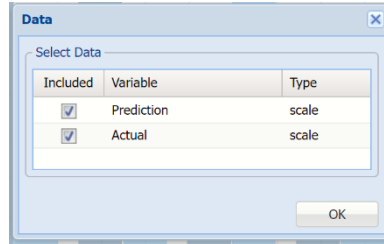
*Bin Details View with continuous (scale) variables*



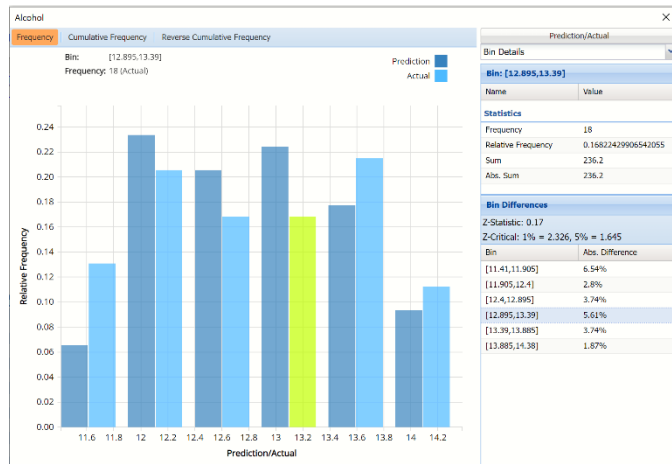
- Frequency is the number of observations assigned to the bin.
- Relative Frequency is the number of observations assigned to the bin divided by the total number of observations.
- Sum is the sum of all observations assigned to the bin.

- Absolute Sum is the sum of the absolute value of all observations assigned to the bin, i.e.  $|observation\ 1| + |observation\ 2| + |observation\ 3| + \dots$

**Bin Differences:** Click *Prediction* to open the Data dialog. Select both Prediction and Actual to add the predicted values for the training partition to the chart.



Two grids are displayed, Bin and Bin Differences. Bin Differences displays the differences between the relative frequencies of each bin for the two histograms, sorted in the same order as the bins listed in the chart. The computed Z-Statistic as well as the critical values, are displayed in the title of the grid.

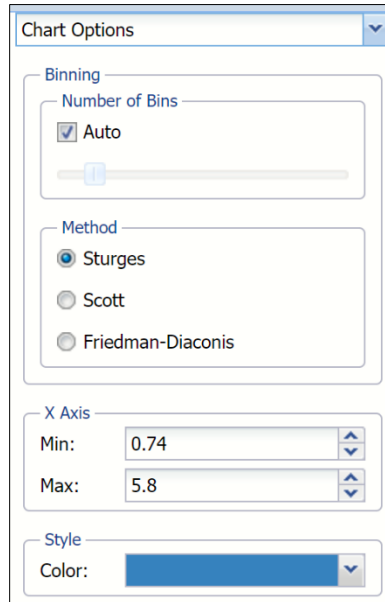


For more information on Bin Details, see the *Generate Data* chapter within the **Data Science Reference Guide**.

### Chart Settings View

The Chart Options view contains controls that allow you to customize the appearance of the charts that appear in the dialog. When you change option selections or type numerical values in these controls, the chart area is instantly updated.

### Chart Options Pane



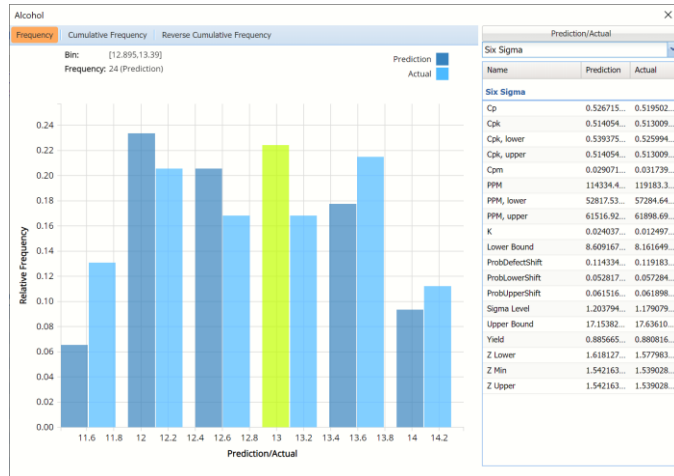
The controls are divided into three groups: Binning, Method and Style.

- **Binning:** Applies to the number of bins in the chart.
  - **Auto:** Select Auto to allow Analytic Solver to automatically select the appropriate number of bins to be included in the frequency charts. See Method below for information on how to change the bin generator used by Analytic Solver when this option is selected.
  - **Manually select # of Bins:** To manually select the number of bins used in the frequency charts, uncheck “Auto” and drag the slider to the right to increase the number of bins or to the left to decrease the number of bins.
- **Method:** Three generators are included in the Analyze Data application to generate the “optimal” number of bins displayed in the chart. All three generators implicitly assume a normal distribution. Sturges is the default setting. The Scott generator should be used with random samples of normally distributed data. The Freedman-Diaconis’ generator is less sensitive than the standard deviation to outliers in the data.
- **X Axis:** Analytic Solver allows users to manually set the Min and Max values for the X Axis. Simply type the desired value into the appropriate text box.
- **Style:**
  - **Color:** Select a color, to apply to the entire variable graph, by clicking the down arrow next to Color and then selecting the desired hue.

Notice in the screenshot below that both the Prediction and Actual data appear in the chart together, and statistics, percentiles or Six Sigma indices for both data appear on the right.

To remove either the Predicted or the Actual data from the chart, click Prediction/Actual in the top right and then uncheck the data type to be removed.

Frequency Chart shown with Prediction and Training data, Six Sigma displayed.



Click the down arrow next to Statistics to view Percentiles or Six Sigma indices for each type of data.

**Training: Prediction Summary and Prediction Details**

The Prediction Summary for the Training Partition lists the following metrics: SSE, MSE, RMSE, MAD and R2. See definitions above.

| Metric | Value    |
|--------|----------|
| SSE    | 4.408646 |
| MSE    | 0.041202 |
| RMSE   | 0.202983 |
| MAD    | 0.156692 |
| R2     | 0.93328  |

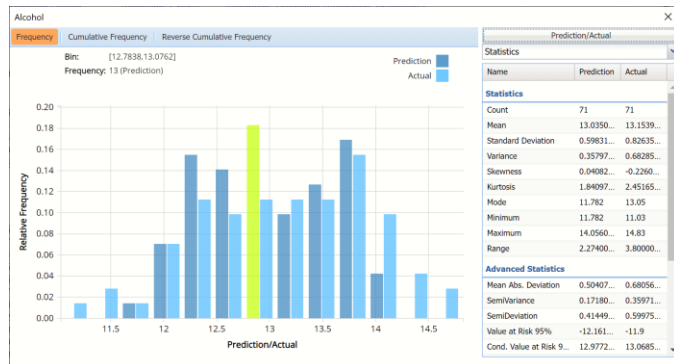
  

| Record ID | Alcohol | Prediction: Alcohol | Residual |
|-----------|---------|---------------------|----------|
| Record 1  | 14.23   | 14.087              | 0.143    |
| Record 5  | 13.24   | 13.277              | -0.037   |
| Record 8  | 14.06   | 13.739              | 0.321    |
| Record 15 | 14.38   | 14.189              | 0.191    |
| Record 18 | 13.83   | 13.744              | 0.086    |

Individual records and their predictions are shown beneath Training: Prediction Details.

**PFBM\_ValidationScore**

**Frequency Chart:** PFBM\_ValidationScore also displays a frequency chart once the tab is selected. See above for an explanation of this chart.



The Prediction Summary for the Validation Partition is shown below.

| Metric | Value    |
|--------|----------|
| SSE    | 26.25256 |
| MSE    | 0.369754 |
| RMSE   | 0.608074 |
| MAD    | 0.480789 |
| R2     | 0.450784 |

| Record ID  | Alcohol | Prediction: Alcohol | Residual |
|------------|---------|---------------------|----------|
| Record 163 | 12.85   | 12.926              | -0.076   |
| Record 126 | 12.07   | 12.468              | -0.398   |
| Record 170 | 13.4    | 13.286              | 0.114    |
| Record 13  | 13.75   | 13.841              | -0.091   |
| Record 39  | 13.05   | 13.442              | -0.392   |

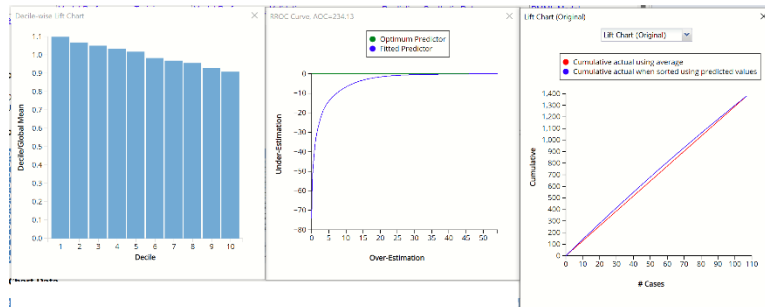
Individual records and their predictions are shown beneath Validation: Prediction Details.

### ***PFBM\_TrainingLiftChart and PFBM\_ValidationLiftChart***

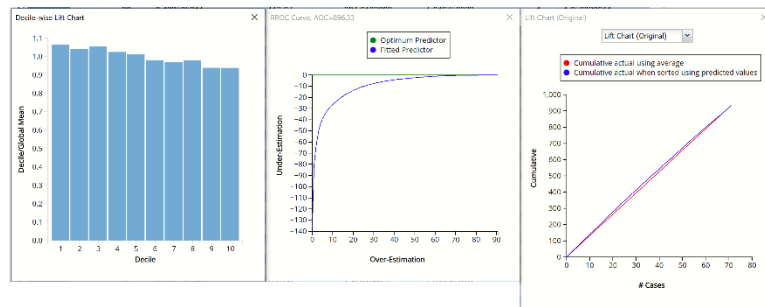
Lift charts and Regression ROC Curves (RROC) are visual aids for measuring model performance. Lift Charts consist of a lift curve and a baseline. The greater the area between the lift curve and the baseline, the better the model. RROC (regression receiver operating characteristic) curves plot the performance of regressors by graphing over-estimations (or predicted values that are too high) versus underestimations (or predicted values that are too low.) The closer the curve is to the top left corner of the graph (in other words, the smaller the area above the curve), the better the performance of the model.

Note: To view these charts in the Cloud app, click the Charts icon on the Ribbon, select PFBM\_TrainingLiftChart or PFBM\_ValidationLiftChart for Worksheet and Decile Chart, RROC Chart or Gain Chart for Chart.

### **Decile-wise Lift Chart, RROC Curve, and Lift Charts for Training Partition**



### **Decile-wise Lift Chart, RROC Curve, and Lift Charts for Valid. Partition**



### **Original Lift Chart**

After the model is built using the training data set, the model is used to score on the training data set and the validation data set (if one exists). Then the data set(s) are sorted using the predicted output variable value. After sorting, the actual outcome values of the output variable are cumulated and the lift curve is drawn as the number of cases versus the cumulated value. The baseline (red line connecting the origin to the end point of the blue line) is drawn as the number of cases versus the average of actual output variable values multiplied by the number of cases.

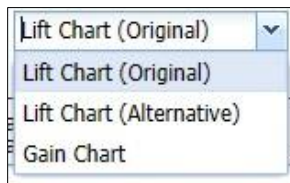
### **Decile-Wise Lift Chart**

The decile-wise lift curve is drawn as the decile number versus the cumulative actual output variable value divided by the decile's mean output variable value. This bars in this chart indicate the factor by which the Linear Regression model outperforms a random assignment, one decile at a time. Typically, this graph will have a "stairstep" appearance - the bars will descend in order from left to right as shown in the decile-wise charts for both partitions.

### **RRoc Curve**

The Regression ROC curve (RROC) was updated in V2017. This new chart compares the performance of the regressor (Fitted Predictor) with an Optimum Predictor Curve. The Optimum Predictor Curve plots a hypothetical model that would provide perfect prediction results. The best possible prediction performance is denoted by a point at the top left of the graph at the intersection of the x and y axis. Area Over the Curve (AOC) is the space in the graph that appears above the RROC curve and is calculated using the formula:  $\sigma^2 * n^2/2$  where n is the number of records. The smaller the AOC, the better the performance of the model.

In V2017, two new charts were introduced: a new Lift Chart and the Gain Chart. To display these new charts, click the down arrow next to Lift Chart (Original), in the Original Lift Chart, then select the desired chart.



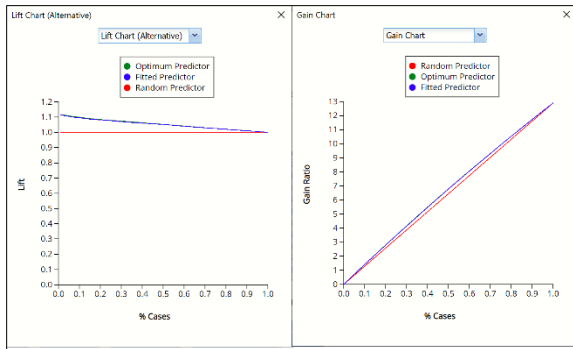
### **Lift Chart Alternative and Gain Chart**

Select Lift Chart (Alternative) to display Analytic Solver Data Science's alternative Lift Chart. Each of these charts consists of an Optimum Predictor curve, a Fitted Predictor curve, and a Random Predictor curve. The Optimum Predictor curve plots a hypothetical model that would provide perfect classification for our data. The Fitted Predictor curve plots the fitted model and the Random Predictor curve plots the results from using no model or by using a random guess (i.e. for x% of selected observations, x% of the total number of positive observations are expected to be correctly classified).

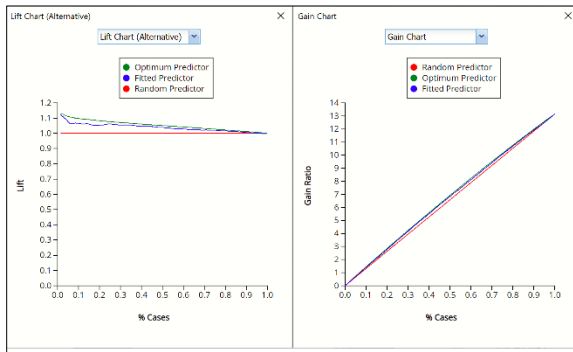
Click the down arrow and select Gain Chart from the menu. In this chart, the Gain Ratio is plotted against the % Cases.



## Lift Chart (Alternative) and Gain Chart for Training Partition



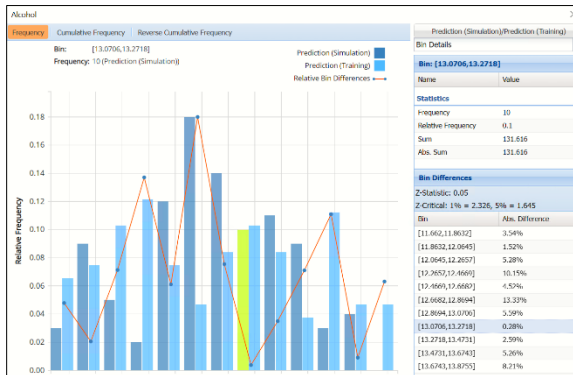
## Lift Chart (Alternative) and Gain Chart for Validation Partition



## PFBM\_Simulation

As discussed above, Analytic Solver Data Science generates a new output worksheet, PFBM\_Simulation, when *Simulate Response Prediction* is selected on the Simulation tab of the Find Best Model dialog.

This report contains the prediction data for the synthetic data, the training data (using the fitted model) and the Excel – calculated Expression column, if populated in the dialog. Users may switch between between the Predicted, Training, and Expression data or a combination of two, as long as they are of the same type. Recall that Expression was not used in this example. For more information on this chart, see above.



Notice that red lines which connect the relative Bin Differences for each bin. Bin Differences are computed based on the frequencies of records which predictions fall into each bin. For example, consider the highlighted bin in the screenshot above  $[x_0, x_1] = [13.0706, 13.2718]$ . There are 10 Simulation

records and 11 Training records in this bin. The relative frequency of the Simulation data is  $10/100 = 10\%$  and the relative frequency of the Training data is  $11/107 = 10.28\%$ . Hence the Absolute Difference (in frequencies) is  $|10 - 10.3| = .28\%$ .

The generated synthetic data is included in the Prediction: Synthetic Data report.

| Record   | Alcohol | Malic_Acid  | Ash      | Ash_Alcalinity | Magnesium    | Total_Phenols | Flavanoids  | Nonflavanoid_Phenols | Proanthocyanins | Color_Intensity | Hue      | OD280_OD315  | Proline      |
|----------|---------|-------------|----------|----------------|--------------|---------------|-------------|----------------------|-----------------|-----------------|----------|--------------|--------------|
| Record 1 | 13.028  | -0.01530916 | -0.35967 | -1.0889645     | 2.12445064   | -0.238975365  | 0.258718473 | -0.78177212          | -0.70663326     | 1.879931243     | 0.440115 | -0.879060321 | 1.307056466  |
| Record 2 | 13.399  | 2.295619103 | -0.32054 | -0.837110946   | 1.109879219  | 0.214484841   | 0.423422026 | -0.653161142         | -0.490766901    | 0.531027284     | -1.02568 | 0.951125374  | 2.487020282  |
| Record 3 | 12.49   | -0.60099888 | -0.06504 | -0.324851639   | -1.479192322 | -1.19268188   | -1.49667675 | 0.22486419           | -1.680420259    | -1.001852639    | -1.1904  | -1.367351721 | -1.163348871 |
| Record 4 | 13.104  | -0.9982244  | -1.30997 | -1.346411042   | 0.63059763   | -0.501791428  | -0.21169036 | -0.344862567         | 0.386979032     | 0.037755079     | 0.029116 | -0.790440567 | -0.362657574 |
| Record 5 | 11.991  | 1.447067647 | 1.316109 | 3.716000049    | 0.48359683   | -0.423390048  | -0.56770029 | 0.297466235          | 0.47699962      | -1.210064041    | 0.115195 | 1.464504456  | -0.989121494 |

## Scoring New Data

Now that the model has been fit to the data, this fitted model will be used to score new patient data, found below. Enter the following new data into a new tab in the workbook.

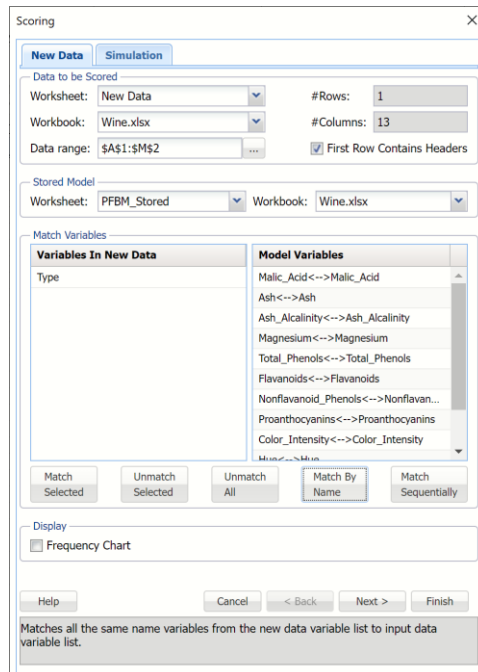
| Type | Malic_Acid | Ash  | Ash_Alcalinity | Magnesium | Total_Phenols | Flavanoids | Nonflavanoid_Phenols | Proanthocyanins | Color_Intensity | Hue  | OD280_OD315 | Proline |
|------|------------|------|----------------|-----------|---------------|------------|----------------------|-----------------|-----------------|------|-------------|---------|
| B    | 2.34       | 2.37 | 19.5           | 100       | 2.30          | 2.03       | 0.36                 | 1.59            | 5.06            | 0.96 | 2.61        | 747     |

Click the New Data tab and then click the Score icon on the Analytic Solver Data Science ribbon.

The screenshot shows the 'Scoring' dialog box with the following callouts:

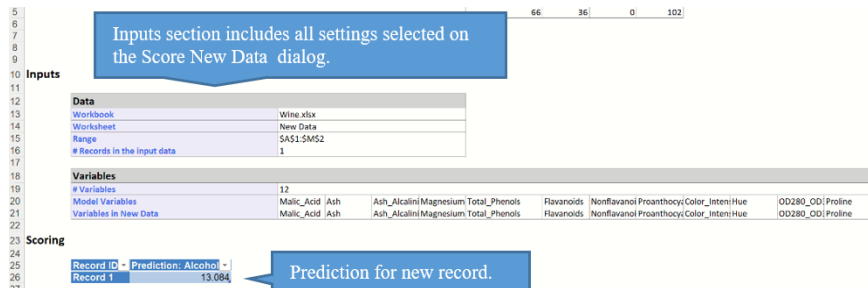
- Name of worksheet containing new data:** Points to the 'New Data' tab.
- Number of rows/columns in new data:** Points to '# Rows: 1' and '# Columns: 13'.
- Name of workbook:** Points to 'Wine.xlsx' in the 'Workbook' field.
- Range for new data:** Points to '\$A\$1:\$M\$2' in the 'Data range' field.
- Variable names in new data:** Points to the 'Variables In New Data' list.
- Variable names in fitted model:** Points to the 'Model Variables' list.

Click "Match By Name" to match each variable in the new data with the same variable in the fitted model, i.e. Malic\_Acid with Malic\_Acid, Ash with Ash, etc.



Click OK to score the new data record and predict the alcohol content of the new wine sample.

A new worksheet, Scoring\_Bagging is inserted to the right.



Notice that the predicted alcohol content for this sample is 13.084.

Please see the “Scoring New Data” chapter within the Analytic Solver Data Science User Guide for information on scoring new data.

# Classifying the Iris Dataset

---

## Introduction

Analytic Solver Data Science's chart wizard allows you to visualize the contents of your dataset by allowing you to create up to 8 different types of charts quickly and easily: bar chart, histogram, line chart, parallel coordinates, scatterplot, scatterplot matrix, boxplot, and variable. Create a bar chart to compare individual statistics (i.e. mean, count, etc.) across a group of variables or a box plot to illustrate the shape of a distribution, its central value and range of data. Use a histogram to depict the range and scale of your observations at variable intervals, a line chart to describe a time series dataset, or a parallel coordinates plot to create a multivariate profile. A scatterplot can be created to compare the relationship between two variables while a scatterplot matrix combines several scatterplots into one panel allowing the user to see pairwise relationships between variables. Finally, use the variable graph to plot each variable's distribution. Two additional options, Export to PowerBI and Export to Tableau, may be used to export your data to these 3<sup>rd</sup> party applications. For more information, see the Analytic Solver Data Science Reference Guide.

The following example will walk you through a quick tutorial, using both Analytic Solver Data Science Cloud and Desktop apps, which create several different types of charts to allow a quick and thorough visualization of the well-known iris dataset.

---

## Creating the Classification Model

In this example, we will use the well-known iris dataset to craft a classification model. The iris dataset is a multivariate data set introduced by Sir Ronald Fisher in 1936 as an example of discriminant analysis. This data set contains 50 observations from each of three Iris species: Iris setosa, Iris virginica and Iris versicolor. Four characteristics were recorded from each sample: sepal length and width and petal length and width. Our goal is to build a model that will assign new data points to the correct iris species.

To open this dataset, click **Help – Example Models – Forecasting/Data Science Examples**, then click the *Iris* link.

A portion of the dataset is shown below.

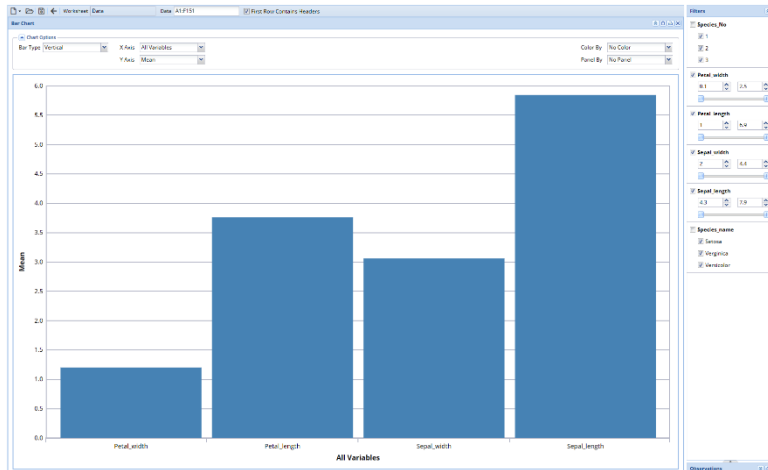
|    | A          | B           | C            | D           | E            | F            |
|----|------------|-------------|--------------|-------------|--------------|--------------|
| 1  | Species_No | Petal_width | Petal_length | Sepal_width | Sepal_length | Species_name |
| 2  | 1          | 0.2         | 1.4          | 3.5         | 5.1          | Setosa       |
| 3  | 1          | 0.2         | 1.4          | 3           | 4.9          | Setosa       |
| 4  | 1          | 0.2         | 1.3          | 3.2         | 4.7          | Setosa       |
| 5  | 1          | 0.2         | 1.5          | 3.1         | 4.6          | Setosa       |
| 6  | 1          | 0.2         | 1.4          | 3.6         | 5            | Setosa       |
| 7  | 1          | 0.4         | 1.7          | 3.9         | 5.4          | Setosa       |
| 8  | 1          | 0.3         | 1.4          | 3.4         | 4.6          | Setosa       |
| 9  | 1          | 0.2         | 1.5          | 3.4         | 5            | Setosa       |
| 10 | 1          | 0.2         | 1.4          | 2.9         | 4.4          | Setosa       |
| 11 | 1          | 0.1         | 1.5          | 3.1         | 4.9          | Setosa       |
| 12 | 1          | 0.2         | 1.5          | 3.7         | 5.4          | Setosa       |
| 13 | 1          | 0.2         | 1.6          | 3.4         | 4.8          | Setosa       |
| 14 | 1          | 0.1         | 1.4          | 3           | 4.8          | Setosa       |
| 15 | 1          | 0.1         | 1.1          | 3           | 4.3          | Setosa       |
| 16 | 1          | 0.2         | 1.2          | 4           | 5.8          | Setosa       |
| 17 | 1          | 0.4         | 1.5          | 4.4         | 5.7          | Setosa       |
| 18 | 1          | 0.4         | 1.3          | 3.9         | 5.4          | Setosa       |
| 19 | 1          | 0.3         | 1.4          | 3.5         | 5.1          | Setosa       |
| 20 | 1          | 0.3         | 1.7          | 3.8         | 5.7          | Setosa       |
| 21 | 1          | 0.3         | 1.5          | 3.8         | 5.1          | Setosa       |
| 22 | 1          | 0.2         | 1.7          | 3.4         | 5.4          | Setosa       |
| 23 | 1          | 0.4         | 1.5          | 3.7         | 5.1          | Setosa       |
| 24 | 1          | 0.2         | 1            | 3.6         | 4.6          | Setosa       |
| 25 | 1          | 0.5         | 1.7          | 3.3         | 5.1          | Setosa       |
| 26 | 1          | 0.2         | 1.9          | 3.4         | 4.8          | Setosa       |
| 27 | 1          | 0.2         | 1.6          | 3           | 5            | Setosa       |
| 28 | 1          | 0.4         | 1.6          | 3.4         | 5            | Setosa       |
| 29 | 1          | 0.2         | 1.5          | 3.5         | 5.2          | Setosa       |
| 30 | 1          | 0.2         | 1.4          | 3.4         | 5.2          | Setosa       |
| 31 | 1          | 0.2         | 1.6          | 3.2         | 4.7          | Setosa       |
| 32 | 1          | 0.2         | 1.6          | 3.1         | 4.8          | Setosa       |
| 33 | 1          | 0.4         | 1.5          | 3.4         | 5.4          | Setosa       |
| 34 | 1          | 0.1         | 1.5          | 4.1         | 5.2          | Setosa       |
| 35 | 1          | 0.2         | 1.4          | 4.2         | 5.5          | Setosa       |

First, we will explore our dataset with graphical analysis. This is often a beginning step when creating a classification or prediction model.

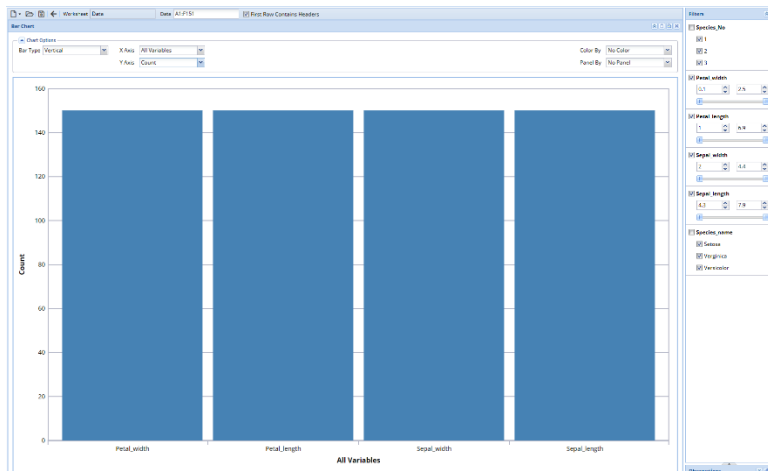
Visualizations can help identify patterns within each variable or in relationships between variables. If a variable is numerical, charts such as histograms or boxplots may be used to study the distribution of values and to identify outliers. If the dataset is a time series, a line chart may be used to detect a trend in the data, for example, the time of year when ski/snowboard sales spike. Bar charts are typically used when a variable is categorical, such as the number of ski's, snowboards, and poles that are purchased during a winter sale. The iris dataset includes both categorical and numerical data. Let's start with a simple bar chart to examine the maximum, minimum and mean values of each numerical variable, petal\_length, petal\_width, sepal\_length and sepal\_width.

- To create a chart, select a cell within the dataset, say cell A2, then click **Explore – Chart Wizard** to open the Chart Wizard. Select **Bar Chart**.
- Uncheck Species\_No and Species\_name from the Filters pane to remove both from the chart.

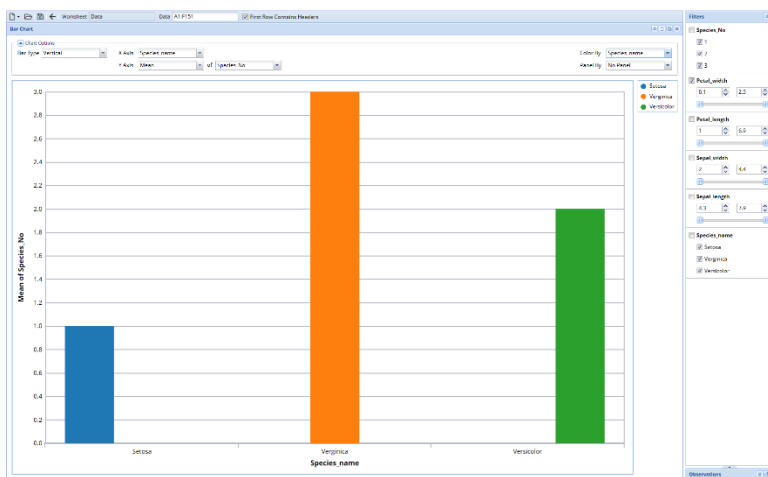
The resultant chart displays the mean values for all variables.



To confirm that all variables have the same data counts, select Count for the Y Axis.



Select Species\_Name for X Axis and Mean of Species\_No for Y Axis. Then set Color By to Species\_name. Then unselect all but Petal\_width in the Filters pane.



This chart becomes a bit more interesting. In this chart, we see that the Setosa species typically has a small petal width, while the Verginica species has a very large petal\_width, in relation to the Setosa species. The petal width of

Versicolor is about midway between the two. This chart supports the idea of including petal\_width as a possible predictor in our classification model.

We could use similar steps to examine the mean values by species of the remaining numerical variables: petal\_length, sepal\_width and sepal\_length.

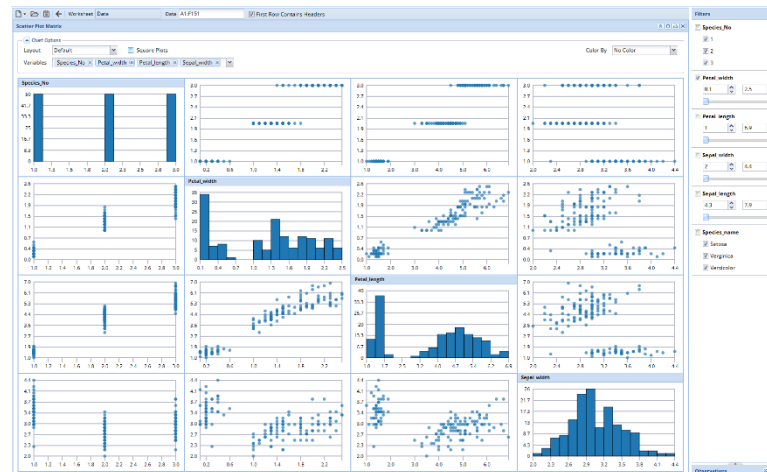
Now let's create a scatterplot matrix to identify any trends in the relationships between variables.

Click the back arrow. You will be asked to Save or Discard the chart. To save, enter a name such as Box Whisker Petal Length in the Name field, then click Save. To discard, click the trash can icon that corresponds to the chart name. To open a previously saved chart, click Explore – Existing Charts.

Select the Scatterpot Matrix.

Uncheck Square Plots at the top to extend the matrix to the width of the chart area.

On the diagonal of the Scatterplot Matrix, we find histograms of each variable. Each histogram displays the frequency of values for each variable.



Now we can clearly see that Setosa iris' have short, narrow petals and short, wide sepals. While the Verginica species has long wide petals and long, medium width sepals.

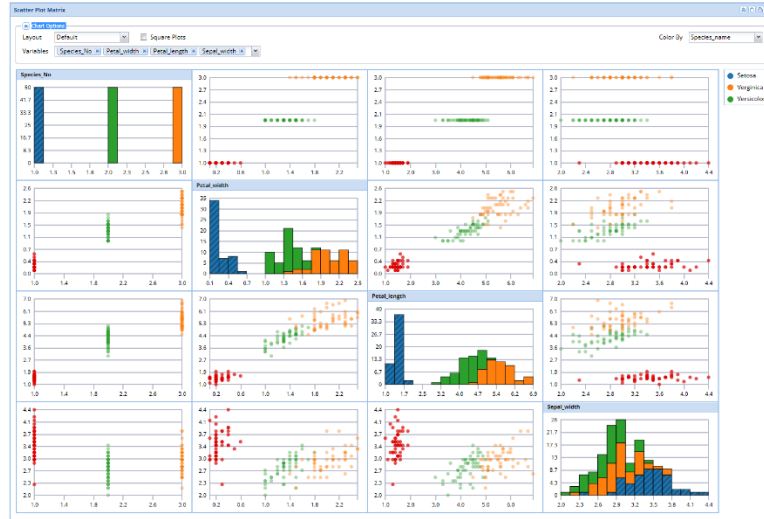
Set Color By to Species Name to distinguish each data point by species.



Moving off to the right of the Petal\_width histogram, we find a scatterplot matrix with Petal\_width on the y axis and petal\_length on the x axis. As you can see there is a distinct cluster of observations with narrow petal\_widths and short petal\_lengths.

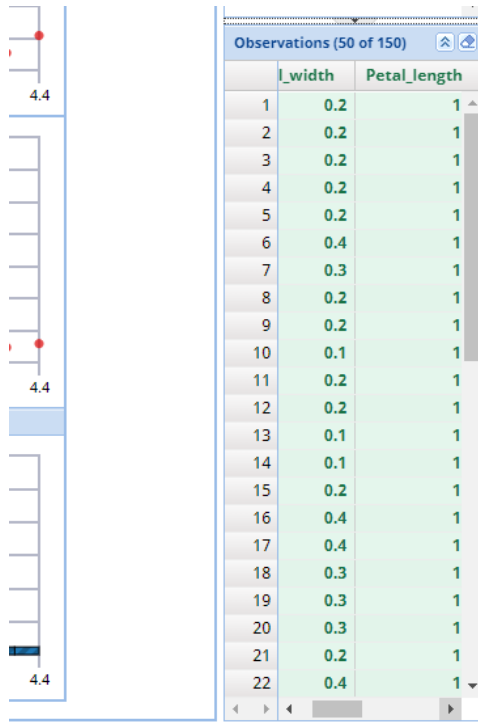
Let's take a closer look by using our mouse to draw a square around the cluster of points in the lower left of this scatterplot in each app.

Immediately, the observations that are in this cluster all turn red, if using the Desktop app, or blue, if using the Cloud app, in each histogram and scatterplot matrix. As you can see, not only do these observations have narrow petal widths and short petal lengths but they also have short sepal lengths with slightly wider sepal widths.

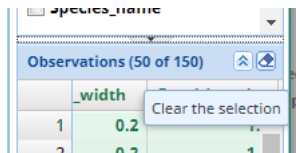


The records for each point in this cluster can be found under Observations on the right of the Chart Wizard. The first thing we see is that this cluster is made up of 50 points out of a total of 150. If we take a moment to scroll through these 50 records using the right and left arrows, one common feature starts to emerge – each of these records belongs to the Setosa species.

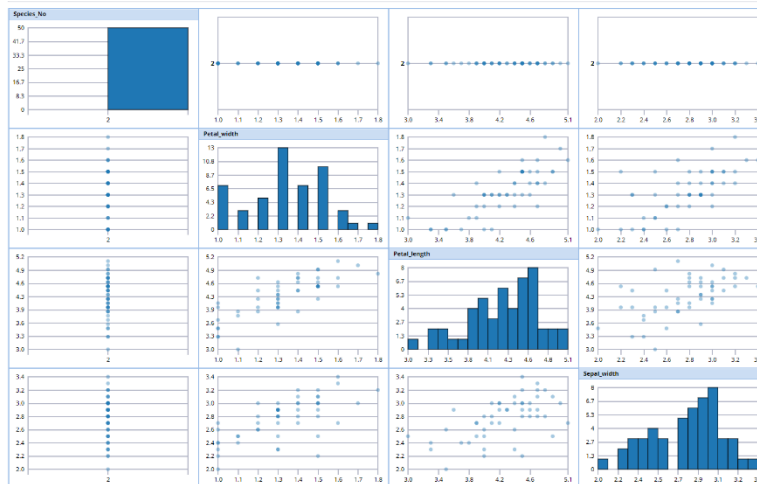




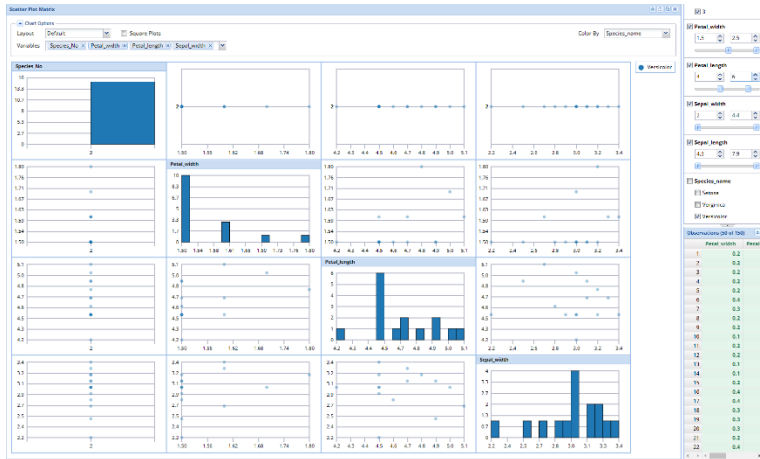
To remove the selection, click the eraser icon in the top, right of the Observations tab.



From here we can start to customize a little deeper in the Destop app. For example, we can look only at the Versicolor species simply by unchecking both Setosa and Verginica under Species\_name.

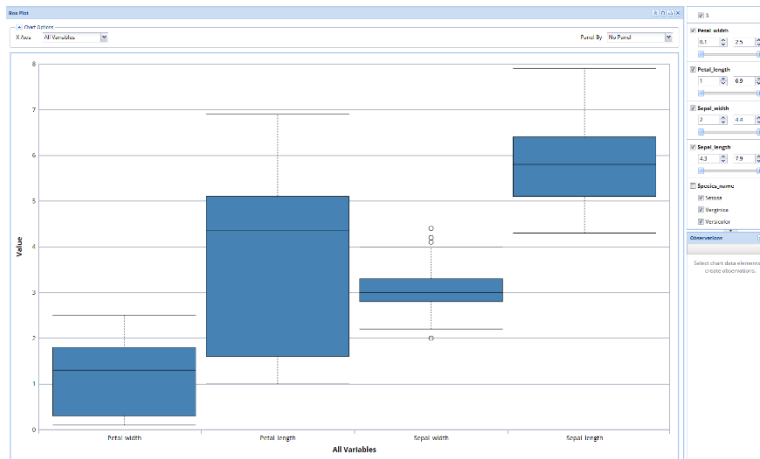


We can further filter our graph by only looking at observations where petal\_width is between 1.5 and 2.5 and petal\_length is between 4 and 6, by either moving the sliders left and right or using the spinner buttons, appropriately.



Add the two species back into the chart and then restore all the filters back to their defaults by simply moving all sliders to their extremes. Now let's create a boxplot, as an alternative to the histograms in the scatterplot matrix. A boxplot is similar to a histogram in that it graphically displays a variable's distribution of values, but a boxplot also includes several additional useful statistics. In a boxplot, also known as a box and whisker plot, the whiskers denote the minimum and maximum values and a "box" is used to designate the 25th and 75th percentiles. The top of the box denotes the 75th percentile and the bottom of the box denotes the 25th percentile. Inside the box is a solid line indicating the variable's median value. Looking at the three boxplots together, we can see that while petal\_widths for Virginica and Versicolor species overlap, while the petal\_widths of the Setosa iris do not.

*Data Science Cloud*



Click the X in the title of the Chart Wizard dialog to close the Chart Wizard. For more information on the chart wizard, including walk throughs on how to create each chart in the wizard, see the Data Science Reference Guide.


In the data science field, datasets with large amounts of variables are routinely encountered. In most cases, the size of the dataset can be reduced by removing highly correlated or superfluous variables. The accuracy and reliability of a classification or prediction model produced from this resultant dataset will be improved by the removal of these redundant and unnecessary variables. In addition, superfluous variables increase the data-collection and data-processing costs of deploying a model on a large database. As a result, one of the first steps in data science should be finding ways to reduce the number of independent or

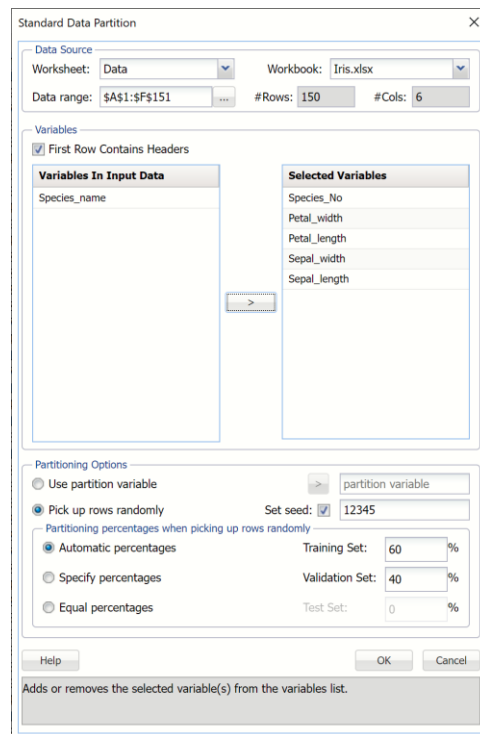
input variables used in the model (otherwise known as dimensionality) without sacrificing accuracy.

Dimensionality Reduction is the process of reducing the amount of variables to be used as input to a prediction or classification model. One way to reduce the number of dimensions in a classification model is to fit a model using a classification tree which splits on the variables that result in the best model. Variables that are not included in the classification tree, can be removed.

One very important issue when fitting a model is how well the newly created model will behave when applied to new data. To address this issue, the dataset is divided into multiple partitions: a training partition which is used to create or "train" the model and a validation partition to test the performance of the model. (For more information on partitioning, see the Analytic Solver Data Science Reference Guide.) For this particular example, we will partition the dataset according to the Standard Partition defaults of 60% of records assigned to the Training set and 40% of records assigned to the Validation set.

Click **Partition – Standard Partition** to open the Standard Data Partition dialog. Select Species\_No, Petal\_width, Petal\_length, Sepal\_width, and

Sepal\_length under Variables In Input data and click  to select and include in the partition. Click OK to accept the defaults and create the partitions.

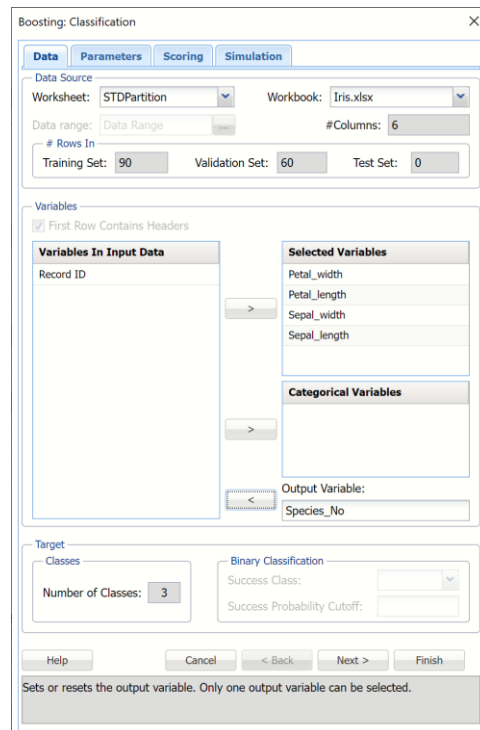


STDPartition is inserted in the Model tab of the Analytic Solver task pane under Data Science – Results – Text Mining. Click the Partition Summary link in the Output Navigator to highlight the records selected for the Training Dataset and click the Validation Data link to highlight the records selected for the Validation Dataset.

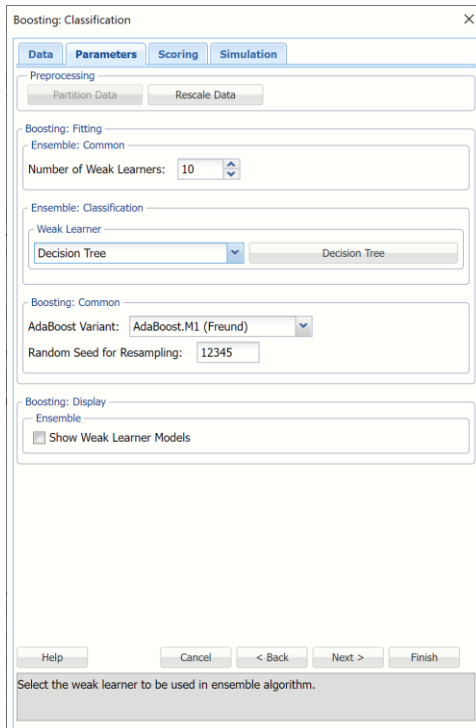


Analytic Solver Data Science offers three powerful ensemble methods for use with Classification trees: bagging (bootstrap aggregating), boosting, and random trees. The Classification Tree Algorithm on its own can be used to find one model that will result in a good classification of the new data. We can view the statistics and confusion matrices of the current classifier to see if our model is a good fit to the data, but how would we know if there is a better classifier? The answer is – we don't. Ensemble methods, however, allow us to combine multiple “weak” classification tree models which, when taken together form a new, more accurate “strong” classification tree model. These methods work by creating multiple diverse classification models, by taking different samples of the original dataset, and then combining their outputs. Outputs may be combined by several techniques, for example, majority vote for classification and averaging for regression. This combination of models effectively reduces the variance in the “strong” model. The three different types of ensemble methods offered in Analytic Solver (bagging, boosting, and random trees) differ on three items: 1. The selection of training data for each classifier or “weak” model, 2. How the “weak” models are generated and 3. How the outputs are combined. In all three methods, each “weak” model is trained on the entire training dataset to become proficient in some portion of the dataset.

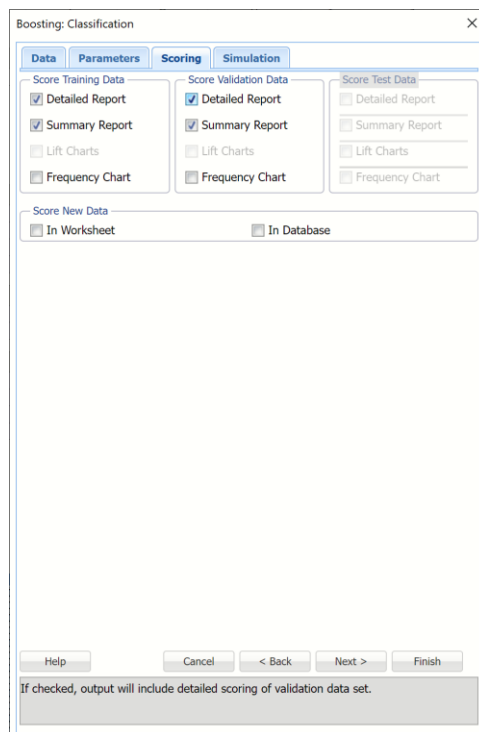
Click Classify – Ensemble – Boosting to open the Boosting – Data tab. Select Petal\_width, Petal\_length, Sepal\_width, and Sepal\_length as Selected Variables and Species\_No as the Output Variable.



Click Next to advance to the Classification Tree Boosting – Parameters tab. Click the down arrow beneath Weak Learner and select Decision Tree from the menu. Note: If we had not previously partitioned the dataset, we could have performed the partitioning from this dialog by clicking Partition Data.



Click Next to accept all default option settings and advance to the Classification Tree Boosting – Scoring tab. Select Detailed Report under both Scoring Training Data and Score Validation Data.



For more information on the Simulation tab and Frequency Chart on the Simulation tab, see the Finding Best Model or the Predicting Housing Prices chapters.

Click Finish to accept all default option settings and run the Boosting algorithm. Output sheets containing results of the algorithm will be inserted to the right of the Data worksheet. Click the Training: Classification Summary and Validation: Classification Summary links in the Output Navigator.

The algorithm performed perfectly on the training data partition as evidenced by the overall error of 0.

| Training: Classification Summary |         |          |         |  |
|----------------------------------|---------|----------|---------|--|
| <b>Confusion Matrix</b>          |         |          |         |  |
| Actual\Predicted                 | 1       | 2        | 3       |  |
| 1                                | 29      | 0        | 0       |  |
| 2                                | 0       | 31       | 0       |  |
| 3                                | 0       | 0        | 30      |  |
| <b>Error Report</b>              |         |          |         |  |
| Class                            | # Cases | # Errors | % Error |  |
| 1                                | 29      | 0        | 0       |  |
| 2                                | 31      | 0        | 0       |  |
| 3                                | 30      | 0        | 0       |  |
| Overall                          | 90      | 0        | 0       |  |
| <b>Metrics</b>                   |         |          |         |  |
| Metric                           | Value   |          |         |  |
| Accuracy (#correct)              | 90      |          |         |  |
| Accuracy (%correct)              | 100     |          |         |  |

In the validation set, four records were assigned incorrectly resulting in a 6.67% misclassification error.

| Validation: Classification Summary |             |          |             |  |
|------------------------------------|-------------|----------|-------------|--|
| <b>Confusion Matrix</b>            |             |          |             |  |
| Actual\Predicted                   | 1           | 2        | 3           |  |
| 1                                  | 21          | 0        | 0           |  |
| 2                                  | 0           | 16       | 3           |  |
| 3                                  | 0           | 1        | 19          |  |
| <b>Error Report</b>                |             |          |             |  |
| Class                              | # Cases     | # Errors | % Error     |  |
| 1                                  | 21          | 0        | 0           |  |
| 2                                  | 19          | 3        | 15.78947368 |  |
| 3                                  | 20          | 1        | 5           |  |
| Overall                            | 60          | 4        | 6.66666667  |  |
| <b>Metrics</b>                     |             |          |             |  |
| Metric                             | Value       |          |             |  |
| Accuracy (#correct)                | 56          |          |             |  |
| Accuracy (%correct)                | 93.33333333 |          |             |  |

If we repeat the same steps using the two remaining Ensemble methods, Bagging (using Decision Trees as the Weak Learner) and Random Trees, we will get the following results.

Each ensemble method performed well, by evidenced by the overall errors in the validation sets, but the random trees method performed best.

### Bagging

| Training: Classification Summary |         |          |         |  |
|----------------------------------|---------|----------|---------|--|
| <b>Confusion Matrix</b>          |         |          |         |  |
| Actual\Predicted                 | 1       | 2        | 3       |  |
| 1                                | 29      | 0        | 0       |  |
| 2                                | 0       | 31       | 0       |  |
| 3                                | 0       | 0        | 30      |  |
| <b>Error Report</b>              |         |          |         |  |
| Class                            | # Cases | # Errors | % Error |  |
| 1                                | 29      | 0        | 0       |  |
| 2                                | 31      | 0        | 0       |  |
| 3                                | 30      | 0        | 0       |  |
| Overall                          | 90      | 0        | 0       |  |
| <b>Metrics</b>                   |         |          |         |  |
| Metric                           | Value   |          |         |  |
| Accuracy (#correct)              | 90      |          |         |  |
| Accuracy (%correct)              | 100     |          |         |  |

| Validation: Classification Summary |             |          |             |  |
|------------------------------------|-------------|----------|-------------|--|
| <b>Confusion Matrix</b>            |             |          |             |  |
| Actual \ Predicted                 | 1           | 2        | 3           |  |
| 1                                  | 21          | 0        | 0           |  |
| 2                                  | 0           | 19       | 0           |  |
| 3                                  | 0           | 2        | 18          |  |
| <b>Error Report</b>                |             |          |             |  |
| Class                              | # Cases     | # Errors | % Error     |  |
| 1                                  | 21          | 0        | 0           |  |
| 2                                  | 19          | 0        | 0           |  |
| 3                                  | 20          | 2        | 10          |  |
| Overall                            | 60          | 2        | 3.333333333 |  |
| <b>Metrics</b>                     |             |          |             |  |
| Metric                             | Value       |          |             |  |
| Accuracy (#correct)                | 58          |          |             |  |
| Accuracy (%correct)                | 96.66666667 |          |             |  |

## Random Trees

| Training: Classification Summary |         |          |         |  |
|----------------------------------|---------|----------|---------|--|
| <b>Confusion Matrix</b>          |         |          |         |  |
| Actual \ Predicted               | 1       | 2        | 3       |  |
| 1                                | 29      | 0        | 0       |  |
| 2                                | 0       | 31       | 0       |  |
| 3                                | 0       | 0        | 30      |  |
| <b>Error Report</b>              |         |          |         |  |
| Class                            | # Cases | # Errors | % Error |  |
| 1                                | 29      | 0        | 0       |  |
| 2                                | 31      | 0        | 0       |  |
| 3                                | 30      | 0        | 0       |  |
| Overall                          | 90      | 0        | 0       |  |
| <b>Metrics</b>                   |         |          |         |  |
| Metric                           | Value   |          |         |  |
| Accuracy (#correct)              | 90      |          |         |  |
| Accuracy (%correct)              | 100     |          |         |  |

| Validation: Classification Summary |         |          |         |  |
|------------------------------------|---------|----------|---------|--|
| <b>Confusion Matrix</b>            |         |          |         |  |
| Actual \ Predicted                 | 1       | 2        | 3       |  |
| 1                                  | 21      | 0        | 0       |  |
| 2                                  | 0       | 19       | 0       |  |
| 3                                  | 0       | 0        | 20      |  |
| <b>Error Report</b>                |         |          |         |  |
| Class                              | # Cases | # Errors | % Error |  |
| 1                                  | 21      | 0        | 0       |  |
| 2                                  | 19      | 0        | 0       |  |
| 3                                  | 20      | 0        | 0       |  |
| Overall                            | 60      | 0        | 0       |  |
| <b>Metrics</b>                     |         |          |         |  |
| Metric                             | Value   |          |         |  |
| Accuracy (#correct)                | 60      |          |         |  |
| Accuracy (%correct)                | 100     |          |         |  |

# Predicting Housing Prices using Multiple Linear Regression

---

## Introduction

Prediction algorithms are supervised learning methods which aim to estimate, or forecast, a continuous output variable. For example, the predicted price of a house that has just come on the market. Analytic Solver Data Science includes 4 different prediction algorithms: Multiple Linear Regression, k-Nearest Neighbors, Regression Tree, and Neural Networks.

This example uses Multiple Linear Regression to fit a prediction model using the Boston\_Housing dataset. The information in this dataset was gathered by the US Census Service from census tracts within the Boston area. Each of the 14 features (or variables) describes a characteristic impacting the selling price of a house. (A description of each variable is given in the example workbook.) In addition to these variables, the data set also contains an additional variable, which has been created by categorizing median value (MEDV) into two categories – high (MEDV > 30) and low (MEDV < 30).

---

## Multiple Linear Regression Example

### Input

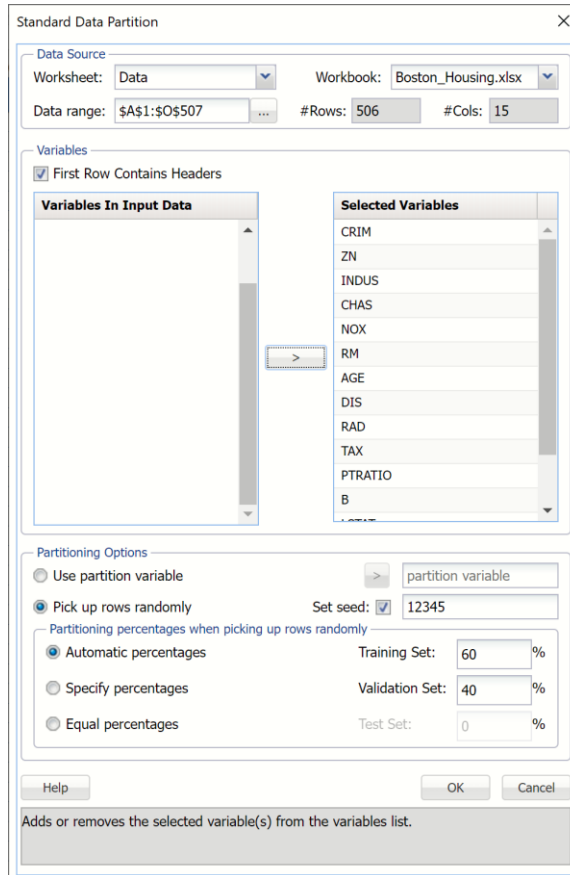
1. Open or upload the Boston\_Housing.xlsx example dataset. A portion of the dataset is shown below. The last variable, CAT.MEDV, is a discrete classification of the MEDV variable and will not be used in this example.

| CRIM    | ZN   | INDUS | CHAS | NOX   | RM    | AGE  | DIS    | RAD | TAX | PTRATIO | B      | LSTAT | MEDV | CAT_MEDV |
|---------|------|-------|------|-------|-------|------|--------|-----|-----|---------|--------|-------|------|----------|
| 0.00632 | 18   | 2.31  | 0    | 0.538 | 6.575 | 65.2 | 4.09   | 1   | 296 | 15.3    | 396.9  | 4.98  | 24   | 0        |
| 0.02731 | 0    | 7.07  | 0    | 0.469 | 6.421 | 78.9 | 4.9671 | 2   | 242 | 17.8    | 396.9  | 9.14  | 21.6 | 0        |
| 0.02729 | 0    | 7.07  | 0    | 0.469 | 7.185 | 61.1 | 4.9671 | 2   | 242 | 17.8    | 392.83 | 4.03  | 34.7 | 1        |
| 0.03237 | 0    | 2.18  | 0    | 0.459 | 6.998 | 45.8 | 6.0622 | 3   | 222 | 18.7    | 394.63 | 2.94  | 33.4 | 1        |
| 0.06905 | 0    | 2.18  | 0    | 0.459 | 7.147 | 54.2 | 6.0622 | 3   | 222 | 18.7    | 396.9  | 5.33  | 36.2 | 1        |
| 0.02985 | 0    | 2.18  | 0    | 0.458 | 6.43  | 58.7 | 6.0622 | 3   | 222 | 18.7    | 384.12 | 5.21  | 28.7 | 0        |
| 0.08829 | 12.5 | 7.87  | 0    | 0.524 | 6.012 | 66.6 | 5.9505 | 5   | 311 | 15.2    | 395.6  | 12.43 | 22.9 | 0        |
| 0.14455 | 12.5 | 7.87  | 0    | 0.524 | 6.172 | 96.1 | 5.9505 | 5   | 311 | 15.2    | 396.9  | 19.15 | 27.1 | 0        |
| 0.21124 | 12.5 | 7.87  | 0    | 0.524 | 5.631 | 100  | 6.0821 | 5   | 311 | 15.2    | 386.63 | 29.93 | 16.5 | 0        |
| 0.17004 | 12.5 | 7.87  | 0    | 0.524 | 6.004 | 85.9 | 6.5921 | 5   | 311 | 15.2    | 386.71 | 17.1  | 18.9 | 0        |
| 0.22489 | 12.5 | 7.87  | 0    | 0.524 | 6.377 | 94.3 | 6.3467 | 5   | 311 | 15.2    | 392.52 | 20.45 | 15   | 0        |
| 0.11747 | 12.5 | 7.87  | 0    | 0.524 | 6.009 | 82.9 | 6.2267 | 5   | 311 | 15.2    | 396.9  | 13.27 | 18.9 | 0        |
| 0.09378 | 12.5 | 7.87  | 0    | 0.524 | 5.889 | 39   | 5.4509 | 5   | 311 | 15.2    | 390.5  | 15.71 | 21.7 | 0        |
| 0.62976 | 0    | 8.14  | 0    | 0.538 | 5.949 | 61.8 | 4.7075 | 4   | 307 | 21      | 396.9  | 8.26  | 20.4 | 0        |
| 0.63796 | 0    | 8.14  | 0    | 0.538 | 6.096 | 84.5 | 4.4619 | 4   | 307 | 21      | 380.02 | 10.26 | 18.2 | 0        |
| 0.62739 | 0    | 8.14  | 0    | 0.538 | 5.834 | 56.5 | 4.4986 | 4   | 307 | 21      | 395.62 | 8.47  | 19.9 | 0        |
| 1.05293 | 0    | 8.14  | 0    | 0.538 | 5.935 | 29.3 | 4.4986 | 4   | 307 | 21      | 386.85 | 6.58  | 23.1 | 0        |
| 0.7842  | 0    | 8.14  | 0    | 0.538 | 5.99  | 61.7 | 4.2579 | 4   | 307 | 21      | 386.75 | 14.67 | 17.5 | 0        |
| 0.80271 | 0    | 8.14  | 0    | 0.538 | 5.456 | 36.6 | 3.7965 | 4   | 307 | 21      | 289.99 | 11.69 | 20.2 | 0        |
| 0.7258  | 0    | 8.14  | 0    | 0.538 | 5.727 | 68.5 | 3.7965 | 4   | 307 | 21      | 380.95 | 11.28 | 18.2 | 0        |
| 1.25179 | 0    | 8.14  | 0    | 0.538 | 5.57  | 98.1 | 3.7979 | 4   | 307 | 21      | 376.57 | 21.02 | 13.6 | 0        |
| 0.85204 | 0    | 8.14  | 0    | 0.538 | 5.965 | 89.2 | 4.0123 | 4   | 307 | 21      | 392.53 | 13.83 | 19.6 | 0        |

2. First, we partition the data into training and validation sets using the Standard Data Partition defaults with percentages of 60% of the data randomly allocated to the Training Set and 40% of the data randomly allocated to the Validation Set. For more information on partitioning a dataset, see the *Data Science Partitioning* chapter within the Analytic Solver Data Science Reference Guide.



### Standard Data Partition Dialog



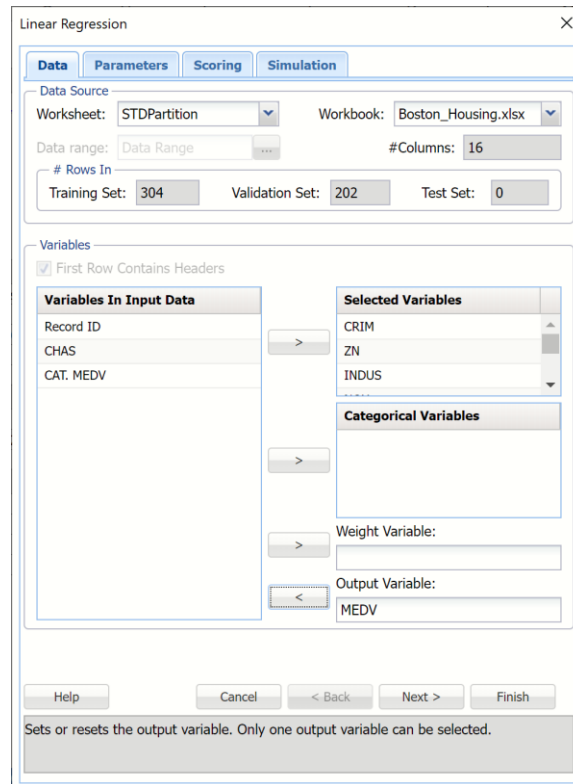
3. With the STDPartition worksheet selected, click **Predict – Linear Regression** to open the Linear Regression Data tab.
4. Select **MEDV** as the *Output Variable*, **CHAS** as a *Categorical Variable* and all remaining variables (except CAT. MEDV) as *Selected variables*.

Note:

All supervised algorithms include a new Simulation tab. This tab uses the functionality from the Generate Data feature (described in the What's New section of this guide and then more in depth in the Analytic Solver Data Science Reference Guide) to generate synthetic data based on the training partition, and uses the fitted model to produce predictions for the synthetic data. The resulting report, LinReq\_Simulation, will contain the synthetic data, the predicted values and the Excel-calculated Expression column, if present. Since this new functionality does not support categorical variables, these types of variables will not be present in the model, only continuous variables.

If the new Simulation feature is not invoked, the categorical variable, CHAS may be added to the model as a categorical variable. However, in this specific instance, since the CHAS variable is the only Categorical Variable selected and its values are binary (0/1), the output results will be the same as if you selected this variable as a Selected (or scale) Variable. If either 1. there were more categorical variables, 2. the CHAS variable was non-binary, or 3, another algorithm besides Multiple Linear Regression or Logistic Regression was selected, this would not be true.

*Linear Regression Dialog, Data tab*



5. Click **Next** to advance to the *Parameters* tab.
6. If *Fit Intercept* is selected, the intercept term will be fitted, otherwise there will be no constant term in the equation. Leave this option selected for this example.
7. Under *Regression: Display*, select all 6 display options to display each in the output.

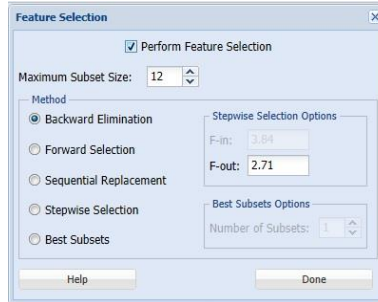
Under *Statistics*, select:

- ANOVA
- Variance-Covariance Matrix
- Multicollinearity Diagnostics

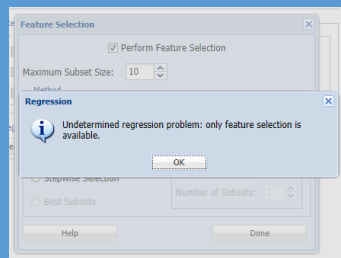
Under *Advanced*, select:

- Analysis of Coefficients
- Analysis of Residuals
- Influence Diagnostics
- Confidence/Prediction Intervals

8. When you have a large number of predictors and you would like to limit the model to only the significant variables, click **Feature Selection** to open the *Feature Selection* dialog and select **Perform Feature Selection** at the top of the dialog. The default setting for *Maximum Subset Size* is **12**. This option can take on values of 1 up to N where N is the number of Selected Variables. The default setting is N.



If the number of rows in the data is less than the number of variables selected as Input variables, the following message box is displayed. Select OK to be directed back to the Feature Selection dialog.

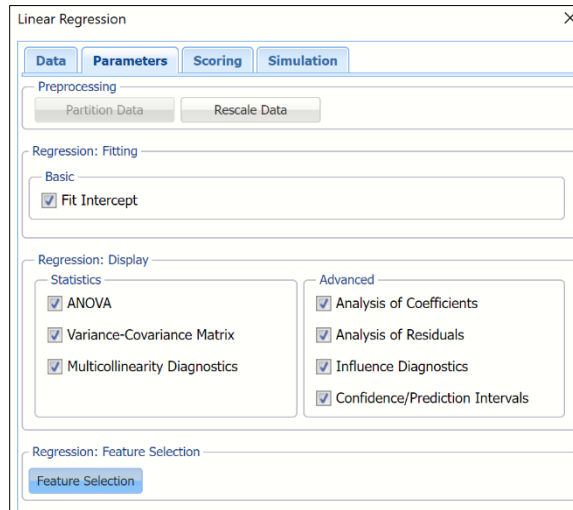


Analytic Solver Data Science offers five different selection procedures for selecting the best subset of variables.

- *Backward Elimination* in which variables are eliminated one at a time, starting with the least significant. If this procedure is selected, FOUT is enabled. A statistic is calculated when variables are eliminated. For a variable to leave the regression, the statistic's value must be less than the value of FOUT (default = 2.71).
- *Forward Selection* in which variables are added one at a time, starting with the most significant. If this procedure is selected, FIN is enabled. On each iteration of the Forward Selection procedure, each variable is examined for the eligibility to enter the model. The significance of variables is measured as a partial F-statistic. Given a model at a current iteration, we perform an F Test, testing the null hypothesis stating that the regression coefficient would be zero if added to the existing set if variables and an alternative hypothesis stating otherwise. Each variable is examined to find the one with the largest partial F-Statistic. The decision rule for adding this variable into a model is: Reject the null hypothesis if the F-Statistic for this variable exceeds the critical value chosen as a threshold for the F Test (FIN value), or Accept the null hypothesis if the F-Statistic for this variable is less than a threshold. If the null hypothesis is rejected, the variable is added to the model and selection continues in the same fashion, otherwise the procedure is terminated.
- *Sequential Replacement* in which variables are sequentially replaced and replacements that improve performance are retained.
- *Stepwise selection* is similar to Forward selection except that at each stage, Analytic Solver Data Science considers dropping variables that are not statistically significant. When this procedure is selected, the Stepwise selection options FIN and FOUT are enabled. In the stepwise selection procedure a statistic is calculated when variables are added or eliminated. For a variable to come into the regression, the statistic's value must be greater than the value for FIN (default = 3.84). For a variable to leave the regression, the statistic's value must be less than the value of FOUT (default = 2.71). The value for FIN must be greater than the value for FOUT.
- *Best Subsets* where searches of all combinations of variables are performed to observe which combination has the best fit. (This option can become quite time consuming depending on the number of input variables.) If this procedure is selected, Number of best subsets is enabled.

9. Click **Done** to accept the default choice, Backward Elimination with a Maximum Subset Size of 3 and an F-out setting of 2.71, and return to the Parameters tab, then click **Next** to advance to the *Scoring* tab.

*Linear Regression Dialog, Parameters tab*



10. Click Next to proceed to the Scoring tab.
11. Select all four options for **Score Training/Validation data**.

When *Detailed report* is selected, Analytic Solver Data Science will create a detailed report of the Discriminant Analysis output.

When *Summary report* is selected (the default), Analytic Solver Data Science will create a report summarizing the Discriminant Analysis output.

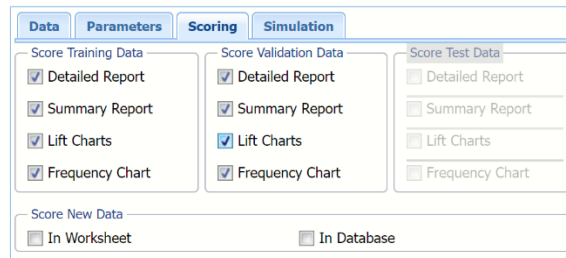
When *Lift Charts* is selected, Analytic Solver Data Science will include Lift Chart and ROC Curve plots in the output.

When Frequency Chart is selected, a frequency chart will be displayed when the LinReg\_TrainingScore and LinReg\_ValidationScore worksheets are selected. This chart will display an interactive application similar to the Analyze Data feature, explained in detail in the Analyze Data chapter that appears earlier in this guide. This chart will include frequency distributions of the actual and predicted responses individually, or side-by-side, depending on the user's preference, as well as basic and advanced statistics for variables, percentiles, six sigma indices.

Since we did not create a test partition, the options for Score test data are disabled. See the chapter "Data Science Partitioning" for information on how to create a test partition.

See the *Scoring New Data* chapter within the Analytic Solver Data Science User Guide for more information on *Score New Data in* options.

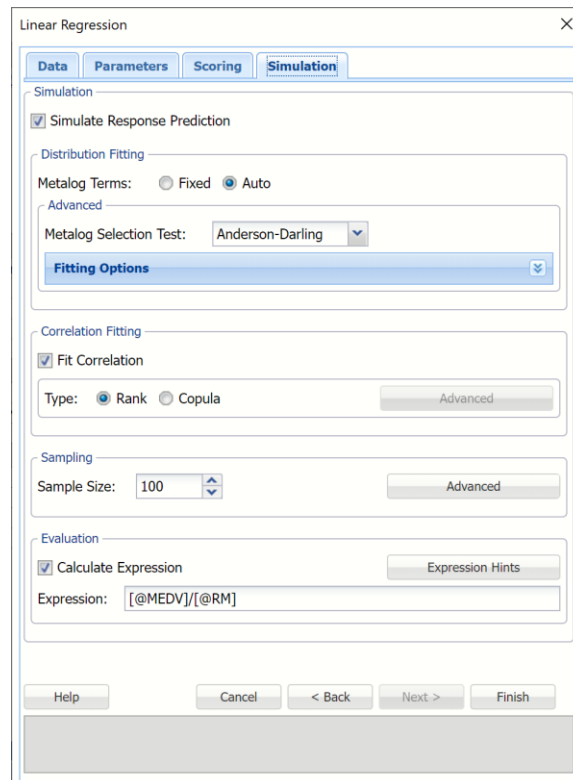
*LinReg dialog, Scoring tab*



12. Click **Next** to advance to the Simulation tab.
13. Select Simulation Response Prediction to enable all options on the Simulation tab of the Linear Regression dialog.

**Simulation tab:** All supervised algorithms include a new Simulation tab. This tab uses the functionality from the Generate Data feature (described earlier in this guide) to generate synthetic data based on the training partition, and uses the fitted model to produce predictions for the synthetic data. The resulting report, LinReg\_Simulation, will contain the synthetic data, the predicted values and the Excel-calculated Expression column, if present. In addition, frequency charts containing the Predicted, Training, and Expression (if present) sources or a combination of any pair may be viewed, if the charts are of the same type.

*Linear Regression dialog, Simulation tab*



**Evaluation:** Select Calculate Expression to amend an Expression column onto the frequency chart displayed on the LinReg\_Simulation output sheet. Expression can be any valid Excel formula that references a variable and the



The R-squared value shown here is the r-squared value for a logistic regression model, defined as -

$$R^2 = (D_0 - D) / D_0,$$

where D is the Deviance based on the fitted model and  $D_0$  is the deviance based on the null model. The null model is defined as the model containing no predictor variables apart from the constant.

|    | B                         | C             | D            |
|----|---------------------------|---------------|--------------|
| 77 | <b>Regression Summary</b> |               |              |
| 78 |                           |               |              |
| 79 |                           | <b>Metric</b> | <b>Value</b> |
| 80 |                           | Residual Df   | 291          |
| 81 |                           | R2            | 0.682503625  |
| 82 |                           | Adjusted R2   | 0.669410991  |
| 83 |                           | Std. Error E  | 5.264721753  |
| 84 |                           | RSS           | 8065.732885  |
| 85 |                           |               |              |

- Predictor Screening:** In Analytic Solver Data Science, a new preprocessing feature selection step was added in V2015 to take advantage of automatic variable screening and elimination, using Rank-Revealing QR Decomposition. This allows Analytic Solver Data Science to identify the variables causing multicollinearity, rank deficiencies, and other problems that would otherwise cause the algorithm to fail. Information about “bad” variables is used in Variable Selection and Multicollinearity Diagnostics, and in computing other reported statistics. Included and excluded predictors are shown in the Model Predictors table. In this model, all predictors were included in the model; all predictors were eligible to enter the model passing the tolerance threshold. This denotes a tolerance beyond which a variance – covariance matrix is not exactly singular to within machine precision. The test is based on the diagonal elements of the triangular factor R resulting from Rank-Revealing QR Decomposition. Predictors that do not pass the test are excluded.

| 86  | <b>Predictor Screening</b>       |             |          |
|-----|----------------------------------|-------------|----------|
| 87  | Predictor                        | Criteria    | Included |
| 88  | Intercept                        | 0.725537578 | TRUE     |
| 89  | CRIM                             | 108.7997552 | TRUE     |
| 90  | ZN                               | 334.3012791 | TRUE     |
| 91  | INDUS                            | 65.43648347 | TRUE     |
| 92  | NOX                              | 1.152404762 | TRUE     |
| 93  | RM                               | 14.56575375 | TRUE     |
| 94  | AGE                              | 447.389624  | TRUE     |
| 95  | DIS                              | 19.19721552 | TRUE     |
| 96  | RAD                              | 57.05078956 | TRUE     |
| 97  | TAX                              | 7790.655556 | TRUE     |
| 98  | PTRATIO                          | 50.06197785 | TRUE     |
| 99  | B                                | 3307.657946 | TRUE     |
| 100 | LSTAT                            | 93.48519887 | TRUE     |
| 101 |                                  |             |          |
| 102 |                                  |             |          |
| 103 |                                  |             |          |
| 104 | Tolerance for entering the model | 5.2587E-10  |          |

Note: If a predictor is excluded, the corresponding coefficient estimates will be 0 in the regression model and the variable – covariance matrix would contain all zeros in the rows and columns that correspond to the excluded predictor. Multicollinearity diagnostics, variable selection and other remaining output will be calculated for the reduced model.

The design matrix may be rank-deficient for several reasons. The most common cause of an ill-conditioned regression problem is the presence of feature(s) that can be exactly or approximately represented by a linear combination of other feature(s). For example, assume that among predictors you have 3 input variables X, Y, and Z where  $Z = a * X + b * Y$  and a and b are constants. This will cause the design matrix to not have a full rank. Therefore, one of these 3 variables will not pass the threshold for entrance and will be excluded from the final regression model.

- **Coefficients:** The Regression Model table contains the estimate, the standard error of the coefficient, the p-value and confidence intervals for each variable included in the model.

| Predictor | Estimate     | Confidence Interval: Lower | Confidence Interval: Upper | Standard Error | T-Statistic  | P-Value      |
|-----------|--------------|----------------------------|----------------------------|----------------|--------------|--------------|
| Intercept | 39.3242992   | 25.04276702                | 53.60575283                | 7.256304724    | 5.419323116  | 1.25647E-07  |
| CRIM      | -0.117543783 | -0.221514987               | -0.0135684579              | 0.052823533    | -2.225317168 | 0.02863303   |
| ZN        | 0.073593787  | 0.032819113                | 0.114368461                | 0.020712764    | 3.552292785  | 0.000445327  |
| INDUS     | 0.089039976  | -0.087142895               | 0.265222847                | 0.089517014    | 0.994670976  | 0.320722678  |
| NOX       | -18.90757499 | -29.61778154               | -8.197368433               | 5.44176459     | -3.474530122 | 0.000589633  |
| RM        | 3.21925831   | 2.01635047                 | 4.422895614                | 0.611371113    | 5.260238089  | 2.71152E-07  |
| AGE       | -0.844862299 | -1.209064099               | -0.04866058                | 0.0448776932   | -1.8545105   | 0.0674432467 |
| DIS       | -1.681960315 | -2.270742537               | -1.093178092               | 0.299155225    | -5.622366503 | 4.41504E-08  |
| RAD       | 0.37469157   | 0.183960218                | 0.570972097                | 0.098317193    | 3.839299575  | 0.00015142   |
| TAX       | -0.017021169 | -0.028165768               | -0.00587657                | 0.005662476    | -3.005958975 | 0.002878217  |
| PTRATIO   | -0.044862299 | -0.089064099               | -0.00066058                | 0.0044876932   | -1.8545105   | 0.0674432467 |
| B         | 0.011459749  | 0.003913741                | 0.019005577                | 0.003834062    | 2.988931309  | 0.003038419  |
| LSTAT     | -0.574639513 | -0.714644547               | -0.434634478               | 0.071135364    | -8.07811309  | 1.76548E-14  |

Note: If a variable has been eliminated by Rank-Revealing QR Decomposition, the variable will appear in red in the Coefficients table with a 0 for Estimate, CI Lower, CI Upper, Standard Error and N/A for T-Statistic and P-Value.

The Standard Error is calculated as each variable is introduced into the model beginning with the constant term and continuing with each variable as it appears in the dataset.

Analytic Solver Data Science produces 95% Confidence Intervals for the estimated values. For a given record, the Confidence Interval gives the mean value estimation with 95% probability. This means that with 95% probability, the regression line will pass through this interval.

- **ANOVA:** The ANOVA table includes the degrees of freedom (DF), sum of squares (SS), mean squares (MS), F-Statistic and P-Value.

| Source     | DF  | SS          | MS          | F-Statistic | P-Value     |
|------------|-----|-------------|-------------|-------------|-------------|
| Regression | 12  | 17338.44024 | 1444.87002  | 52.12882473 | 2.88268E-85 |
| Error      | 291 | 8065.732885 | 27.71729514 | N/A         | N/A         |
| Total      | 303 | 25404.17313 | 83.84215553 | N/A         | N/A         |

- **Variance-Covariance Matrix of Coefficients:** Entries in the matrix are the covariances between the indicated coefficients. The “on-diagonal” values are the estimated variances of the corresponding coefficients.

| Predictor | Intercept    | CRIM        | ZN           | INDUS        | NOX          | RM           | AGE         | DIS         | RAD          | TAX          | PTRATIO      | B            | LSTAT        |
|-----------|--------------|-------------|--------------|--------------|--------------|--------------|-------------|-------------|--------------|--------------|--------------|--------------|--------------|
| Intercept | 52.85395824  | 0.00581158  | 0.004181195  | 0.023984102  | -21.45429586 | -1.180309948 | 0.012470154 | -0.20948195 | 0.199903943  | -0.005810274 | -0.043007272 | -0.0006348   | -0.151186028 |
| CRIM      | 0.00581158   | 0.00029059  | 0.00029059   | 0.00029059   | 0.00029059   | 0.00029059   | 0.00029059  | 0.00029059  | 0.00029059   | 0.00029059   | 0.00029059   | 0.00029059   | 0.00029059   |
| ZN        | 0.004181195  | 0.000429205 | 0.000429205  | 0.000429205  | 0.000429205  | 0.000429205  | 0.000429205 | 0.000429205 | 0.000429205  | 0.000429205  | 0.000429205  | 0.000429205  | 0.000429205  |
| INDUS     | 0.023984102  | 0.00031367  | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367  | 0.00031367  | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367   |
| NOX       | -21.45429586 | 0.000211933 | 0.000211933  | 0.000211933  | 0.000211933  | 0.000211933  | 0.000211933 | 0.000211933 | 0.000211933  | 0.000211933  | 0.000211933  | 0.000211933  | 0.000211933  |
| RM        | -1.180309948 | 0.001311193 | 0.001311193  | 0.001311193  | 0.001311193  | 0.001311193  | 0.001311193 | 0.001311193 | 0.001311193  | 0.001311193  | 0.001311193  | 0.001311193  | 0.001311193  |
| AGE       | 0.012470154  | 0.000239841 | 0.000239841  | 0.000239841  | 0.000239841  | 0.000239841  | 0.000239841 | 0.000239841 | 0.000239841  | 0.000239841  | 0.000239841  | 0.000239841  | 0.000239841  |
| DIS       | -0.20948195  | 0.001563441 | 0.001563441  | 0.001563441  | 0.001563441  | 0.001563441  | 0.001563441 | 0.001563441 | 0.001563441  | 0.001563441  | 0.001563441  | 0.001563441  | 0.001563441  |
| RAD       | 0.199903943  | 0.00031367  | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367  | 0.00031367  | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367   | 0.00031367   |
| TAX       | -0.005810274 | 4.99213E-05 | -3.25594E-05 | -0.000219061 | -0.000960202 | 0.000281951  | 1.26031E-06 | 0.00120488  | -0.000439687 | 3.20035E-05  | -0.000117458 | 5.26291E-07  | -1.36477E-07 |
| PTRATIO   | -0.043007272 | 5.3698E-05  | 0.000101778  | 0.001462342  | 0.014969032  | 0.017228223  | -0.00089794 | 0.000191199 | -0.00217458  | 0.00017458   | 0.00017458   | -3.90102E-06 | -0.00081925  |
| B         | -0.0006348   | 7.02724E-05 | -6.1127E-07  | 1.50308E-05  | 0.000337637  | 1.7299E-06   | 5.22891E-05 | 2.84892E-05 | 5.26291E-07  | -3.98102E-05 | 4.476E-05    | 6.0599E-05   | 0.0006348    |
| LSTAT     | -0.151186028 | 0.000960208 | 0.000960208  | 0.000960208  | 0.000960208  | 0.000960208  | 0.000960208 | 0.000960208 | 0.000960208  | 0.000960208  | 0.000960208  | 0.000960208  | 0.000960208  |

- **Multicollinearity Diagnostics:** This table helps assess whether two or more variables so closely track one another as to provide essentially the same information.

| Row ID | Component #1     | Component #2 | Component #3 | Component #4 | Component #5 | Component #6 | Component #7 | Component #8 | Component #9  | Component #10 | Component #11 | Component #12 | Component #13 | Component #14 |
|--------|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|---------------|---------------|---------------|---------------|---------------|---------------|
| 133    | Component #1     | Component #2 | Component #3 | Component #4 | Component #5 | Component #6 | Component #7 | Component #8 | Component #9  | Component #10 | Component #11 | Component #12 | Component #13 | Component #14 |
| 134    | Condition Number | 88.81577303  | 37.73109382  | 30.02073837  | 28.76949098  | 19.75820929  | 15.80446419  | 11.95414593  | 9.874817981   | 8.02777719    | 6.534823285   | 5.84422759    | 3.127125149   | 2.5187983     |
| 135    | Intercept        | 0.00143834   | 0.001623148  | 0.0017912    | 0.001991134  | 0.00226037   | 0.00262541   | 0.003027399  | 0.00361042    | 0.004396313   | 0.00549457    | 0.00684678    | 0.01057106    | 0.01705106    |
| 136    | CRIM             | 0.00106384   | 0.00194353   | 0.002727     | 0.00359873   | 0.00459389   | 0.00576759   | 0.00716255   | 0.00884791    | 0.01081409    | 0.01317573    | 0.01594485    | 0.02061617    | 0.02835412    |
| 137    | ZN               | 0.00258237   | 0.002997544  | 0.00352899   | 0.00419837   | 0.00491452   | 0.00578712   | 0.006852881  | 0.008170755   | 0.009787938   | 0.0117411893  | 0.000202589   | 0.001958248   | 0.00034373    |
| 138    | INDUS            | 0.00078488   | 0.002303317  | 0.00321744   | 0.004237675  | 0.005389141  | 0.006701804  | 0.008258837  | 0.010047554   | 0.012092623   | 0.014464E-08  | 0.000132328   | 1.48811E-05   | 7.47901E-07   |
| 139    | NOX              | 0.55433877   | 0.348021915  | 0.046709976  | 0.00193104   | 0.03871419   | 0.007915296  | 0.000727416  | 0.002072712   | 0.0009982     | 0.000202089   | 5.02151E-05   | 7.38233E-06   | 7.8959E-05    |
| 140    | RM               | 0.021079445  | 0.000496032  | 0.003097912  | 0.004761245  | 0.00682333   | 0.00941882   | 0.012385758  | 0.016020228   | 0.000405546   | 0.000202958   | 0.000223573   | 0.000202729   | 0.000405595   |
| 141    | AGE              | 0.003794538  | 0.003471179  | 0.004967465  | 0.006880188  | 0.009487382  | 0.01226982   | 0.016237037  | 0.020848097   | 0.026471213   | 0.033684608   | 0.04268E-05   | 3.24245E-05   | 0.00010277    |
| 142    | DIS              | 0.08277785   | 0.00380839   | 0.288187109  | 0.36423798   | 0.524217278  | 0.009808502  | 0.049384911  | 0.039677748   | 0.057414821   | 0.002464176   | 0.000420474   | 0.000425429   | 0.003127351   |
| 143    | RAD              | 0.017248031  | 0.065249822  | 0.45207485   | 0.20651153   | 0.008871959  | 0.000411445  | 0.000261973  | 0.0004049E-05 | 0.001198709   | 0.004892296   | 0.000969E-05  | 1.82344E-05   | 0.00010401    |
| 144    | TAX              | 0.34402779   | 0.36037439   | 0.144809883  | 0.44869185   | 0.004432657  | 0.00041277   | 1.7489E-05   | 0.00140234    | 0.00071659    | 0.00022381    | 0.000101331   | 1.9277E-05    | 2.8207E-05    |
| 145    | PTRATIO          | 0.072217634  | 0.017448884  | 0.008923795  | 0.027898805  | 0.04628515   | 0.017188335  | 0.007486976  | 0.024818443   | 0.009533322   | 0.008828203   | 0.001330685   | 0.00017893    | 0.00101968    |
| 146    | B                | 0.00379357   | 0.219523232  | 0.025889154  | 0.052985916  | 0.052875255  | 0.007907514  | 0.0128591    | 0.024869429   | 0.004178934   | 0.004795466   | 0.000460223   | 2.81017E-06   | 0.01137891    |
| 147    | LSTAT            | 0.00204222   | 2.78844E-06  | 0.001209921  | 4.50985E-05  | 0.005132907  | 6.7071E-06   | 0.0001028    | 3.9880E-05    | 0.0001028     | 3.9880E-05    | 6.9369E-05    | 1.7759E-05    | 1.6466E-05    |
| 148    | CHAS_0           | 0.98855244   | 1.99302E-05  | 0.000503742  | 0.001192994  | 0.006733447  | 0.000275304  | 1.13957E-05  | 0.000450615   | 0.00041432    | 0.000152803   | 0.00088788    | 0.020114391   | 7.47148E-06   |
| 149    |                  |              |              |              |              |              |              |              |               |               |               |               |               | 1.56626E-05   |

The columns represent the variance components (related to principal components in multivariate analysis), while the rows represent the variance proportion decomposition explained by each variable in the model. The eigenvalues are those associated with the singular value decomposition of the variance-covariance matrix of the coefficients, while the condition numbers are the ratios of the square root of the largest eigenvalue to all the



rest. In general, multicollinearity is likely to be a problem with a high condition number (more than 20 or 30), and high variance decomposition proportions (say more than 0.5) for two or more variables.

### LinReg\_FS

Select the **Feature Selection** link on the Output Navigator to display the Variable Selection table. This table displays a list of different models generated using the selections made on the *Feature Selection* dialog. When *Backward elimination* is used, Linear Regression may stop early when there is no variable eligible for elimination as evidenced in the table below (i.e. there are no subsets with less than 12 coefficients). Since *Fit Intercept* was selected on the *Parameters* tab, each subset includes an intercept.

| Feature Selection |           |      |    |       |     |    |     |     |     |     |         |   |       |        |        |
|-------------------|-----------|------|----|-------|-----|----|-----|-----|-----|-----|---------|---|-------|--------|--------|
| Best Subsets      |           |      |    |       |     |    |     |     |     |     |         |   |       |        |        |
| Subset ID         | Intercept | CRIM | ZN | INDUS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | CHAS_0 | CHAS_1 |
| Subset 1          | 0         | 1    | 1  | 1     | 1   | 1  | 1   | 0   | 1   | 1   | 1       | 1 | 1     | 1      | 1      |
| Subset 2          | 0         | 1    | 1  | 0     | 1   | 1  | 1   | 0   | 1   | 1   | 1       | 1 | 1     | 1      | 1      |

| Best Subsets Details |               |          |              |                |                         |             |
|----------------------|---------------|----------|--------------|----------------|-------------------------|-------------|
| Subset ID            | #Coefficients | RSS      | Mallows's Cp | R <sup>2</sup> | Adjusted R <sup>2</sup> | Probability |
| Subset 1             | 13            | 7794.742 | 12.03227708  | 0.693171       | 0.680518074             | 0.857546341 |
| Subset 2             | 12            | 7801.43  | 10.28115486  | 0.692908       | 0.68133998              | 0.868915558 |

The error values for the Best Subsets are:

- **RSS:** The residual sum of squares, or the sum of squared deviations between the predicted probability of success and the actual value (1 or 0)
- **Cp:** Mallows Cp (Total squared error) is a measure of the error in the best subset model, relative to the error incorporating all variables. Adequate models are those for which Cp is roughly equal to the number of parameters in the model (including the constant), and/or Cp is at a minimum
- **R-Squared:** R-squared Goodness-of-fit
- **Adj. R-Squared:** Adjusted R-Squared values.
- **"Probability"** is a quasi hypothesis test of the proposition that a given subset is acceptable; if Probability < .05 we can rule out that subset.

Compare the RSS value as the number of coefficients in the subset increases from 13 to 12 (7794.742 to 7801.43). The RSS for 12 coefficients is just slightly higher than the RSS for 13 coefficients suggesting that a model with 12 coefficients may be sufficient to fit a regression.

### LinReg\_ResidInfluence

This output sheet includes two output tables: Residuals and Influence Diagnostics.

- **Residuals:** Click the **Residuals** link in the Output Navigator to open the Residuals table. This table displays the Raw Residuals, Standardized Residuals, Studentized Residuals and Deleted Residuals.

| Residuals |              |              |              |            |  |
|-----------|--------------|--------------|--------------|------------|--|
| Record ID | Raw          | Standardized | Studentized  | Deleted    |  |
| Record 1  | -5.583599752 | -1.060568823 | -1.075954954 | -1.0762476 |  |
| Record 5  | 9.068656439  | 1.722532902  | 1.747874192  | 1.75410039 |  |
| Record 8  | 8.165012831  | 1.550891617  | 1.592572068  | 1.59680726 |  |
| Record 11 | -3.107791911 | -0.590305064 | -0.60907608  | -0.6084166 |  |
| Record 12 | 2.22877225   | 0.441405527  | 0.45895810   | 0.4592400  |  |

Studentized residuals are computed by dividing the unstandardized residuals by quantities related to the diagonal elements of the hat matrix, using a common scale estimate computed without the  $i^{\text{th}}$  case in the model. These residuals have  $t$  - distributions with  $(n-k-1)$  degrees of freedom. As a result, any residual with absolute value exceeding 3 usually requires attention.

The Deleted residual is computed for the  $i^{\text{th}}$  observation by first fitting a model without the  $i^{\text{th}}$  observation, then using this model to predict the  $i^{\text{th}}$  observation. Afterwards the difference is taken between the predicted observation and the actual observation.

- **Influence Diagnostics:** Click the Influence Diagnostics link on the Output Navigator to display the Influence Diagnostics data table. This table contains various statistics computed by Analytic Solver Data Science.

The Cooks Distance for each observation is displayed in this table. This is an overall measure of the impact of the  $i^{\text{th}}$  datapoint on the estimated regression coefficient. In linear models Cooks Distance has, approximately, an  $F$  distribution with  $k$  and  $(n-k)$  degrees of freedom.

|     | B                            | C         | D               | E            | F                | G          | H                 |
|-----|------------------------------|-----------|-----------------|--------------|------------------|------------|-------------------|
| 318 | <b>Influence Diagnostics</b> |           |                 |              |                  |            |                   |
| 319 |                              |           |                 |              |                  |            |                   |
| 320 |                              |           |                 |              |                  |            |                   |
| 321 |                              | Record ID | Cook's Distance | DFFITS       | Covariance Ratio | Leverage   | Delete-1 Variance |
| 322 |                              | Record 1  | 0.002602581     | -0.183989023 | 1.02197407       | 0.02839546 | 27.70222466       |
| 323 |                              | Record 5  | 0.00696548      | 0.301989254  | 0.938715735      | 0.0287865  | 27.52087851       |
| 324 |                              | Record 8  | 0.010627541     | 0.372684627  | 0.984117158      | 0.0516586  | 27.57046173       |
| 325 |                              | Record 11 | 0.001843706     | -0.154648972 | 1.095022209      | 0.06068787 | 27.77741551       |
| 326 |                              | Record 12 | 0.000676929     | -0.093680138 | 1.081173689      | 0.04149535 | 27.79344374       |

The DF fits for each observation is displayed in the output. DFFits gives information on how the fitted model would change if a point was not included in the model.

Analytic Solver Data Science computes DFFits using the following computation.

$$DFFits_i = \frac{\hat{y}_i - \hat{y}_i(-i)}{\sigma(-i)\sqrt{h_i}} = \dots = \frac{\sqrt{h_i}e_i}{\sigma(-i)(1-h_i)} = \dots = e_i^{stud} \sqrt{\frac{h_i}{1-h_i}}$$

where

$\hat{y}_i$  =  $i$ -th fitted value from full model

$\hat{y}_i(-i)$  =  $i$ -th fitted value from model not including  $i$ -th observation

$\sigma(-i)$  = estimated error variance of model not including  $i$ -th observation

$h_i$  = leverage of  $i$ -th point (i.e.  $\{i,i\}$ -th element of Hat Matrix)

$e_i$  =  $i$ -th residual from full model

$e_i^{stud}$  =  $i$ -th Studentized residual

The covariance ratios are displayed in the output. This measure reflects the change in the variance-covariance matrix of the estimated coefficients when the  $i^{\text{th}}$  observation is deleted.

The diagonal elements of the hat matrix are displayed under the Leverage column. This measure is also known as the Leverage of the  $i^{\text{th}}$  observation.

### LinReg Intervals

Click either the **Intervals: Training** or **Intervals: Validation** links in the Output Navigator to view the Intervals report for both the Training and Validation partitions. Of primary interest in a data-science context will be the

predicted and actual values for each record, along with the residual (difference) and Confidence and Prediction Intervals for each predicted value.

| Intervals: Training |                       |                       |                       |                       |  |
|---------------------|-----------------------|-----------------------|-----------------------|-----------------------|--|
| Record ID           | 95% Confidence: Lower | 95% Confidence: Upper | 95% Prediction: Lower | 95% Prediction: Upper |  |
| Record 1            | 27.6247766            | 31.07580706           | 19.00200099           | 39.69838267           |  |
| Record 5            | 25.27808265           | 28.743607             | 16.66142645           | 37.36026321           |  |
| Record 8            | 16.10782046           | 20.78519178           | 7.978581466           | 28.91443078           |  |
| Record 11           | 15.13555428           | 20.19262256           | 7.152119964           | 28.17605688           |  |
| Record 12           | 18.64106281           | 22.84076716           | 10.32376568           | 31.15809429           |  |

| Intervals: Validation |                       |                       |                       |                       |  |
|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|--|
| Record ID             | 95% Confidence: Lower | 95% Confidence: Upper | 95% Prediction: Lower | 95% Prediction: Upper |  |
| Record 222            | 32.01150152           | 36.90807056           | 23.96683159           | 44.95274049           |  |
| Record 104            | 18.96136819           | 22.09719074           | 10.20618198           | 30.85237696           |  |
| Record 163            | 37.47107593           | 43.82141999           | 29.96028926           | 51.33220667           |  |
| Record 411            | 10.11691424           | 20.79715244           | 3.940746268           | 26.97332042           |  |
| Record 460            | 17.21298656           | 20.45374923           | 8.502170883           | 29.16456649           |  |

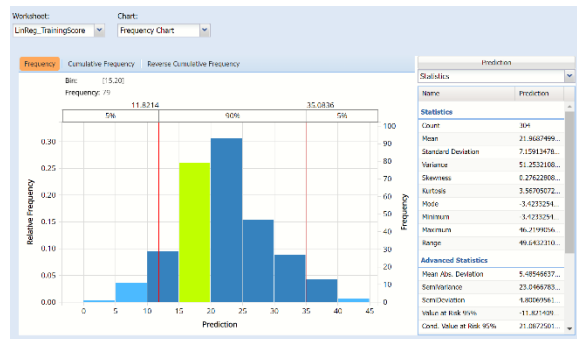
Analytic Solver Data Science produces 95% Confidence and Prediction Intervals for the predicted values. Typically, Prediction Intervals are more widely utilized as they are a more robust range for the predicted value. For a given record, the Confidence Interval gives the mean value estimation with 95% probability. This means that with 95% probability, the regression line will pass through this interval. The Prediction Interval takes into account possible future deviations of the predicted response from the mean. There is a 95% chance that the predicted value will lie within the Prediction interval.

### LinReg\_TrainingScore

Of primary interest in a data-science context will be the predicted and actual values for the MEDV variable along with the residual (difference) for each predicted value in the Training partition.

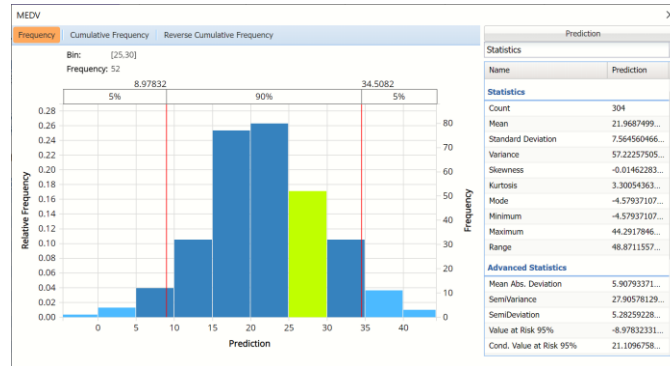
*LinReg\_TrainingScore* displays the newly added Output Variable frequency chart, the Training: Prediction Summary and the Training: Prediction Details report. All calculations, charts and predictions on the *LinReg\_TrainingScore* output sheet apply to the Training partition.

Note: To view charts in the Cloud app, click the Charts icon on the Ribbon, select a worksheet under Worksheet and a chart under Chart.



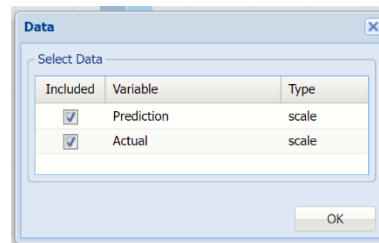
- Frequency Charts:** The output variable frequency chart for the training partition opens automatically once the *LinReg\_TrainingScore* worksheet is selected. To close this chart, click the “x” in the upper right hand corner of the chart. To reopen, click onto another tab and then click back to the *LinReg\_TrainingScore* tab. To move the dialog to a new location on the screen, simply grab the title bar and drag the dialog to the desired location. This chart displays a detailed, interactive frequency chart for the predicted values in the simulated, or synthetic data, and for the predicted values in the training partition.

*Training (Actual) data, MEDV variable*

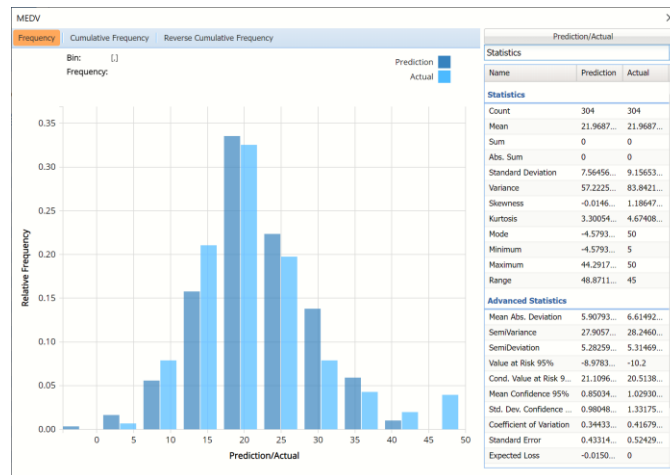


To display the predicted values in the synthetic and training data, click Prediction in the upper right hand corner and select both checkboxes in the Data dialog.

*Click Actual to add Prediction data to the interactive chart*



*Actual vs Predicted for Training Partition, MEDV variable*



Notice in the screenshot below that both the Actual and Prediction data appear in the chart together, and statistics for both data appear on the right.

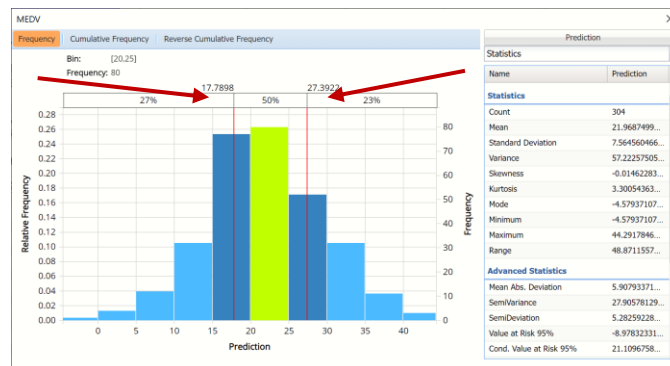
| Prediction/Actual          |            |            |
|----------------------------|------------|------------|
| Statistics                 |            |            |
| Name                       | Prediction | Actual     |
| <b>Statistics</b>          |            |            |
| Count                      | 304        | 304        |
| Mean                       | 21.9687... | 21.9687... |
| Standard Deviation         | 7.56456... | 9.15653... |
| Variance                   | 57.2225... | 83.8421... |
| Skewness                   | -0.0146... | 1.18647... |
| Kurtosis                   | 3.30054... | 4.67408... |
| Mode                       | -4.5793... | 50         |
| Minimum                    | -4.5793... | 5          |
| Maximum                    | 44.2917... | 50         |
| Range                      | 48.8711... | 45         |
| <b>Advanced Statistics</b> |            |            |
| Mean Abs. Deviation        | 5.90793... | 6.61492... |
| SemiVariance               | 27.9057... | 28.2460... |
| SemiDeviation              | 5.28259... | 5.31469... |

To remove either the Actual or Prediction data from the chart, click Prediction/Actual in the top right and then uncheck the data type to be removed.

This chart behaves the same as the interactive chart in the Analyze Data feature found on the Explore menu (and explained in depth in the Data Science Reference Guide).

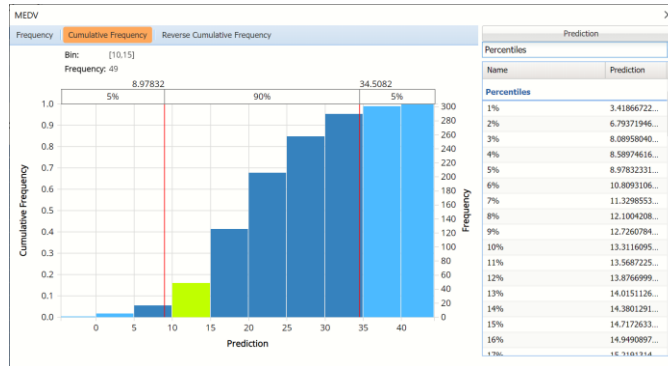
- Use the mouse to hover over any of the bars in the graph to populate the Bin and Frequency headings at the top of the chart.
- When displaying either Original or Synthetic data (not both), red vertical lines will appear at the 5% and 95% percentile values in all three charts (Frequency, Cumulative Frequency and Reverse Cumulative Frequency) effectively displaying the 90<sup>th</sup> confidence interval. The middle percentage is the percentage of all the variable values that lie within the ‘included’ area, i.e. the darker shaded area. The two percentages on each end are the percentage of all variable values that lie outside of the ‘included’ area or the “tails”. i.e. the lighter shaded area. Percentile values can be altered by moving either red vertical line to the left or right.

*Frequency Chart for MEDV Prediction with red percentile lines moved*



- Click Cumulative Frequency and Reverse Cumulative Frequency tabs to see the Cumulative Frequency and Reverse Cumulative Frequency charts, respectively.

*Cumulative Frequency chart with Percentiles displayed*



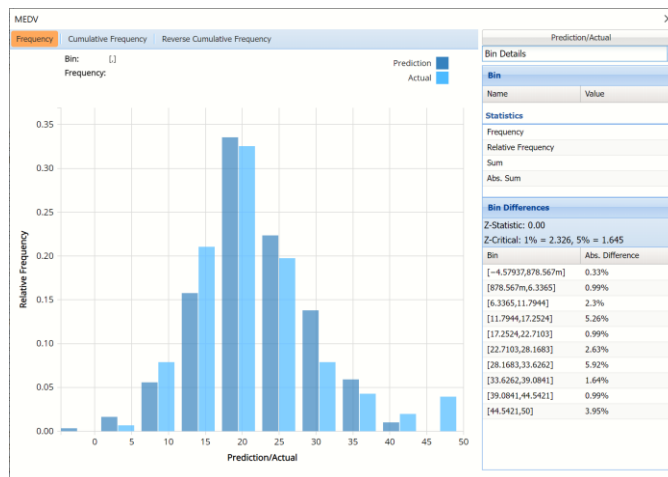
- Click the down arrow next to Statistics to view Percentiles for each type of data along with Six Sigma indices. Use the Chart Options view to manually select the number of bins to use in the chart, as well as to set personalization options.

*Reverse Cumulative Frequency chart with Six Sigma indices displayed*



- Click the down arrow next to Statistics to view Bin Details for each bin in the chart. If both datasets are included in the chart, the Bin Differences are added to the pane.

*Frequency chart with Bin Details displayed*



As discussed above, see the [Analyze Data](#) section of the Exploring Data chapter in the Data Science Reference Guide for an in-depth discussion of

this chart as well as descriptions of all statistics, percentiles, bin details and six sigma indices.

- **Prediction Summary:** In the Prediction Summary report, Analytic Solver Data Science displays the total sum of squared errors summaries for the training partition. The total sum of squared errors is the sum of the squared errors (deviations between predicted and actual values) and the root mean square error (square root of the average squared error). The average error is typically very small, because positive prediction errors tend to be counterbalanced by negative ones.

*Training Prediction Summary*

|    | B                                   | C             | D            |
|----|-------------------------------------|---------------|--------------|
| 9  |                                     |               |              |
| 10 | <b>Training: Prediction Summary</b> |               |              |
| 11 |                                     |               |              |
| 12 |                                     | <b>Metric</b> | <b>Value</b> |
| 13 |                                     | SSE           | 8065.733     |
| 14 |                                     | MSE           | 26.53202     |
| 15 |                                     | RMSE          | 5.150924     |
| 16 |                                     | MAD           | 3.522646     |
| 17 |                                     | R2            | 0.682504     |

- **Prediction Details:** Scroll down to the Training: Prediction Details report to find the Prediction value for the MEDV variable for each record, as well as the Residual value.

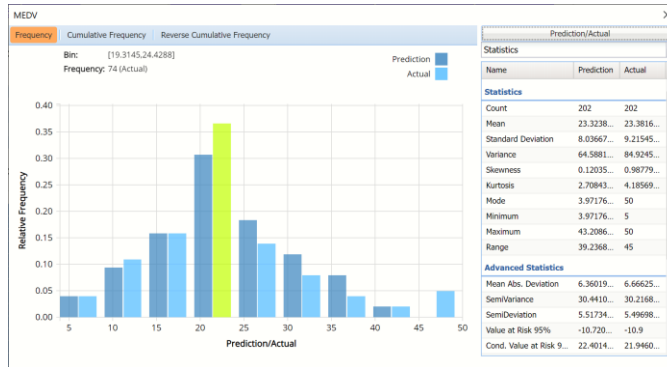
|    | <b>Training: Prediction Details</b> |      |                  |            |
|----|-------------------------------------|------|------------------|------------|
|    | Record ID                           | MEDV | Prediction: MEDV | Residual   |
| 22 | Record 1                            | 24   | 29.58359975      | -5.5835998 |
| 23 | Record 5                            | 36.2 | 27.13134356      | 9.0686564  |
| 24 | Record 8                            | 27.1 | 18.93498717      | 8.1650128  |
| 25 | Record 11                           | 15   | 18.10779191      | -3.1077919 |

### ***LinReg\_ValidationLiftChart***

Another interest in a data-science context will be the predicted and actual values for the MEDV variable along with the residual (difference) for each predicted value in the Validation partition.

*LinReg\_ValidationScore* displays the newly added Output Variable frequency chart, the Training: Prediction Summary and the Training: Prediction Details report. All calculations, charts and predictions on the *LinReg\_ValidationScore* output sheet apply to the Validation partition.

- **Frequency Charts:** The output variable frequency chart for the validation partition opens automatically once the *LinReg\_ValidationScore* worksheet is selected. This chart displays a detailed, interactive frequency chart for the Actual variable data and the Predicted data, for the validation partition. For more information on this chart, see the *LinReg\_TrainingLiftChart* explanation above.



- **Prediction Summary:** In the Prediction Summary report, Analytic Solver Data Science displays the total sum of squared errors summaries for the Validation partition.

| Metric | Value    |
|--------|----------|
| SSE    | 3413.538 |
| MSE    | 16.8987  |
| RMSE   | 4.110803 |
| MAD    | 3.070292 |
| R2     | 0.800025 |

- **Prediction Details:** Scroll down to the Training: Prediction Details report to find the Prediction value for the MEDV variable for each record, as well as the Residual value.

| Record ID  | MEDV | Prediction: MEDV | Residual   |
|------------|------|------------------|------------|
| Record 229 | 46.7 | 34.68191295      | 12.018087  |
| Record 104 | 19.3 | 20.60066938      | -1.3006694 |
| Record 163 | 50   | 37.84226115      | 12.157739  |
| Record 411 | 15   | 15.4593185       | -0.4593185 |
| Record 460 | 20   | 19.0848984       | 0.9151016  |

### ***LinReg\_TrainingLiftChart & LinReg\_ValidationLiftChart***

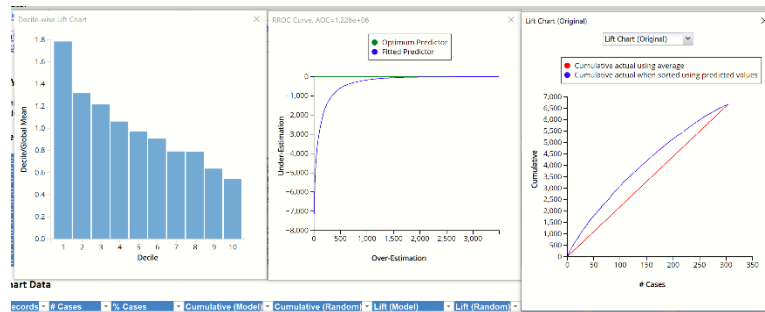
Lift charts and RROC Curves (on the *LinReg\_TrainingLiftChart* and *LinReg\_ValidationLiftChart* tabs, respectively) are visual aids for measuring model performance. Lift Charts consist of a lift curve and a baseline. The greater the area between the lift curve and the baseline, the better the model. RROC (regression receiver operating characteristic) curves plot the performance of regressors by graphing over-estimations (or predicted values that are too high) versus underestimations (or predicted values that are too low.) The closer the curve is to the top left corner of the graph (in other words, the smaller the area above the curve), the better the performance of the model.

After the model is built using the training data set, the model is used to score on the training data set and the validation data set (if one exists). Then the data set(s) are sorted using the predicted output variable value. After sorting, the actual outcome values of the output variable are cumulated and the lift curve is drawn as the number of cases versus the cumulated value. The baseline (red line connecting the origin to the end point of the blue line) is drawn as the number of cases versus the average of actual output variable values multiplied by the number of cases. The decilewise lift curve is drawn as the decile number versus

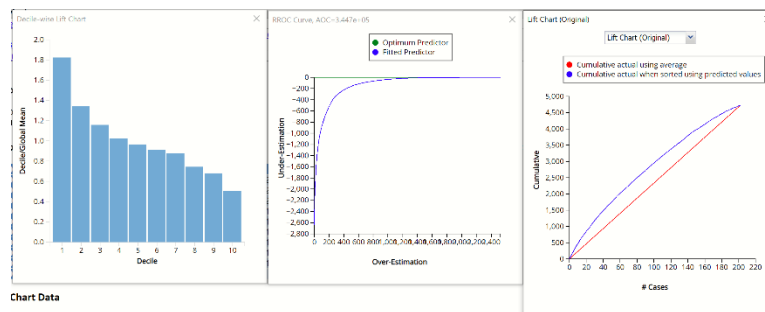


the cumulative actual output variable value divided by the decile's mean output variable value. The bars in this chart indicate the factor by which the MLR model outperforms a random assignment, one decile at a time. Refer to the validation graph below. In the first decile in the validation dataset, taking the most expensive predicted housing prices in the dataset, the predictive performance of the model is about 1.8 times better as simply assigning a random predicted value.

### Decile-Wise Lift Chart, ROC Curve and Lift Chart from Training Partition



### Decile-Wise Lift Chart, ROC Curve and Lift Chart from Validation Partition



In an RROC curve, we can compare the performance of a regressor with that of a random guess (red line) for which under estimations are equal to over-estimations shifted to the minimum under estimate. Anything to the left of this line signifies a better prediction and anything to the right signifies a worse prediction. The best possible prediction performance would be denoted by a point at the top left of the graph at the intersection of the x and y axis. This point is sometimes referred to as the “perfect classification”. Area Over the Curve (AOC) is the space in the graph that appears above the ROC curve and is calculated using the formula:  $\sigma^2 * n^2 / 2$  where  $n$  is the number of records. The smaller the AOC, the better the performance of the model. In this example we see that the area above the curve in both datasets, or the AOC, is fairly large which indicates that this model might not be the best fit to the data.

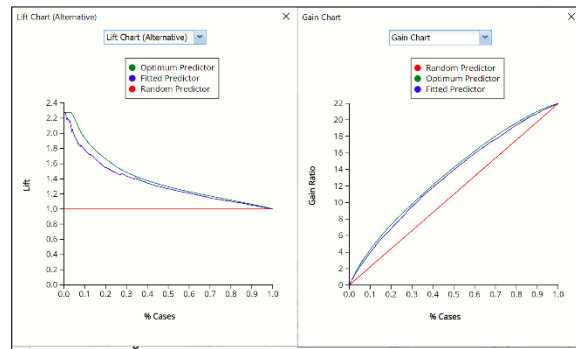
Two new charts were introduced in V2017: a new Lift Chart and the Gain Chart. To display these new charts, click the down arrow next to Lift Chart (Original), in the Original Lift Chart, then select the desired chart.



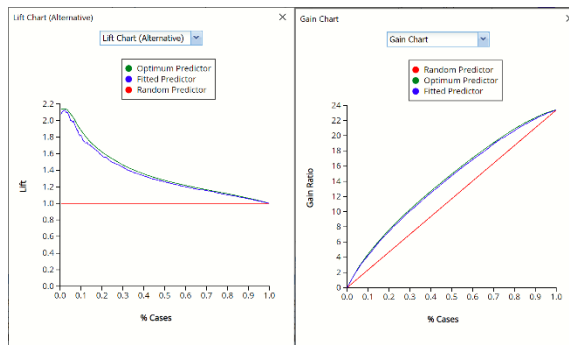
Select Lift Chart to display Analytic Solver Data Science's new Lift Chart. Each of these charts consists of an Optimum Predictor curve, a Fitted Predictor curve, and a Random Predictor curve. The Optimum Predictor curve plots a hypothetical model that would provide perfect classification for our data. The Fitted Predictor curve plots the fitted model and the Random Predictor curve plots the results from using no model or by using a random guess (i.e. for x% of selected observations, x% of the total number of positive observations are expected to be correctly classified).

The Alternative Lift Chart plots Lift against % Cases. The Gain chart plots Gain Ratio against % Cases.

### Lift Chart (Alternative) and Gain Chart for Training Partition



### Lift Chart (Alternative) and Gain Chart for Validation Partition



### LinReg\_Simulation

As discussed above, Analytic Solver Data Science generates a new output worksheet, LinReg\_Simulation, when *Simulate Response Prediction* is selected on the Simulation tab of the Linear Regression dialog.

This report contains the synthetic data, the predicted values for the simulation, or synthetic, data, the predicted values for the training partition (using the fitted model) and the Excel – calculated Expression column, if populated in the dialog. A dialog may be opened to switch between the synthetic data, training data and the expression results, or a combination of two, as long as they are of the same type.

## Synthetic Data

Prediction: Synthetic Data

| Record    | Expression  | MEDV     | CRIM     | ZN       | INDUS    | NOX      | RM       | AGE      | DIS      | RAD      | TAX      | PTRATIO    | B           | LSTAT    |
|-----------|-------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|------------|-------------|----------|
| Record 1  | AGE>=50     | -1.19102 | 0.500923 | 3.99E-06 | 12.18571 | 0.50452  | 4.277229 | 56.19205 | 3.70118  | 1.485154 | 597.5868 | 21.9554523 | 260.6748362 | 25.78433 |
| Record 2  | AGE>=50     | 5.852657 | 1.212722 | 5.67E-06 | 10.78745 | 0.662264 | 5.516223 | 97.77404 | 4.767676 | 16.54911 | 365.3312 | 21.9986003 | 396.8859965 | 31.79133 |
| Record 3  | 40.63775593 | 40.63776 | 0.079214 | 94.88604 | 3.288879 | 0.411282 | 7.59911  | 6.788945 | 9.223803 | 23.99008 | 283.1248 | 16.9199892 | 395.9935904 | 2.891447 |
| Record 4  | 21.8759396  | 21.87594 | 0.568544 | 8.14E-05 | 16.99077 | 0.461248 | 5.231984 | 48.57698 | 4.966019 | 23.98736 | 691.4589 | 18.7598277 | 386.5463592 | 8.549528 |
| Record 5  | AGE>=50     | 33.20467 | 0.019491 | 7.74E-10 | 7.158165 | 0.681577 | 8.132518 | 96.43415 | 2.235082 | 1.022976 | 208.2954 | 17.177104  | 396.8962923 | 6.790191 |
| Record 6  | AGE>=50     | 12.30632 | 37.22072 | 3.02E-09 | 24.41323 | 0.838997 | 7.156659 | 100      | 1.538889 | 23.99968 | 628.0738 | 20.1141971 | 396.6722574 | 28.04089 |
| Record 7  | 36.4147536  | 36.41475 | 0.20491  | 1.35E-05 | 5.231616 | 0.428268 | 6.449745 | 45.69951 | 3.809381 | 21.62634 | 402.0899 | 15.7638283 | 396.8999999 | 4.386995 |
| Record 8  | 24.60900468 | 24.609   | 0.01338  | 0.000549 | 2.528619 | 0.400401 | 6.367766 | 4.391138 | 9.422203 | 1.000269 | 215.1339 | 13.6803202 | 396.631062  | 3.062426 |
| Record 9  | AGE>=50     | 18.03762 | 0.053639 | 0.000291 | 4.395518 | 0.413542 | 5.880412 | 68.4885  | 10.18877 | 3.98318  | 303.5161 | 15.0228847 | 396.8999915 | 7.622623 |
| Record 10 | AGE>=50     | -10.7255 | 0.718453 | 2.38E-11 | 16.30373 | 0.701278 | 3.756476 | 72.15391 | 3.31626  | 1.242232 | 649.1871 | 21.768587  | 358.8771142 | 35.73517 |

Note the first column in the output, Expression. This column was inserted into the Synthetic Data results because Calculate Expression was selected and an Excel function was entered into the Expression field, on the Simulation tab of the Linear Regression dialog

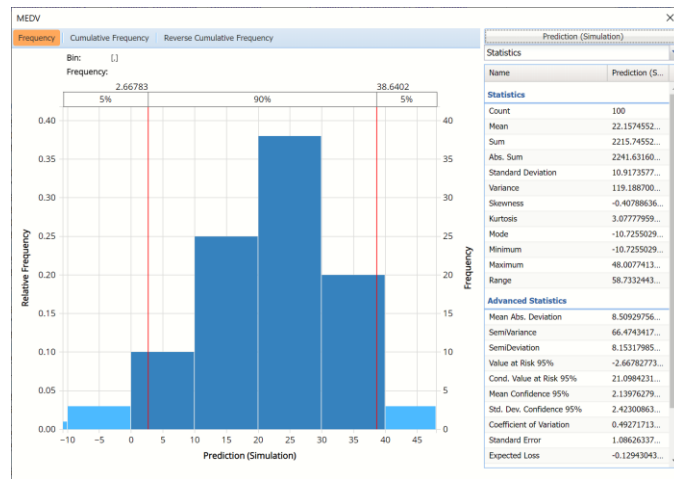
$$[@MEDV]/[@RM]$$

As a result, the values in this column will be equal to the MEDV variable value divided by the number of rooms.

The remainder of the data in this report is synthetic data, generated using the Generate Data feature described in the chapter with the same name, that appears in the Data Science Reference Guide.

The chart that is displayed once this tab is selected, contains frequency information pertaining to the predicted value for the MEDV variable in the synthetic data.

Frequency Chart for Prediction (Simulation) data



Click *Prediction (Simulation)* to add the Actual data to the chart.

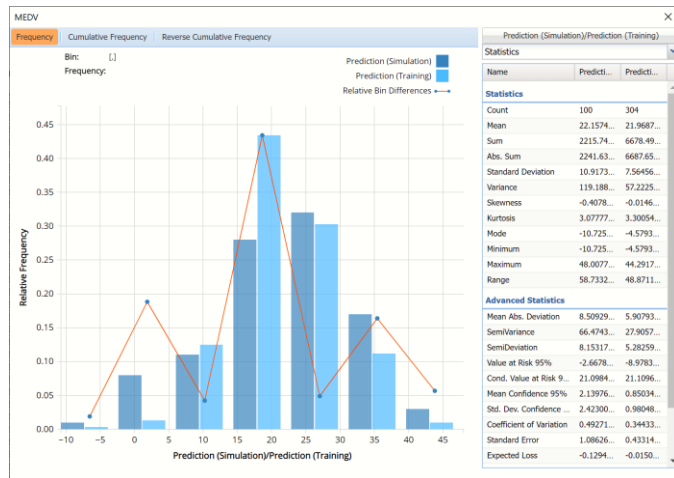
Data dialog

Data

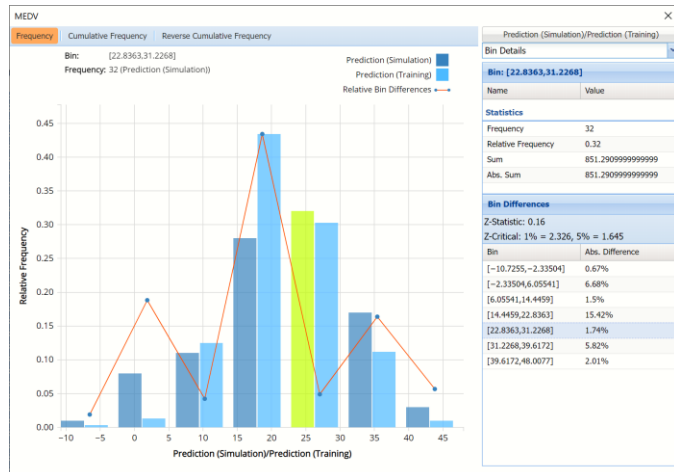
Select Data

| Included                            | Variable                | Type  |
|-------------------------------------|-------------------------|-------|
| <input checked="" type="checkbox"/> | Prediction (Simulation) | scale |
| <input checked="" type="checkbox"/> | Prediction (Training)   | scale |
| <input type="checkbox"/>            | Expression (Simulation) | scale |
| <input type="checkbox"/>            | Expression (Training)   | scale |

Prediction (Simulation) and Prediction (Training) for MEDV variable



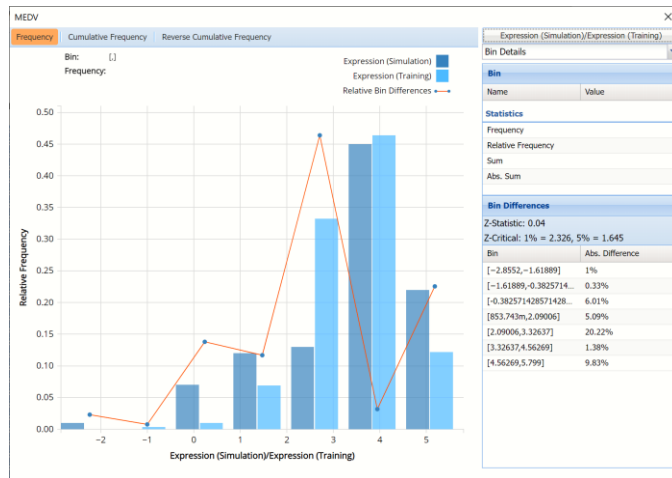
The Relative Bin Differences curve gives the absolute difference between the bins for each dataset. (Click the down arrow next to Frequency and select Bin Details to view.)



Click the down arrow next to Frequency to change the chart view to Relative Frequency or to change the look by clicking Chart Options. Statistics on the right of the chart dialog are discussed earlier in this section. For more information on the generated synthetic data, see the Generate Data chapter that appears earlier in this guide.

Click the Prediction (Simulation)/Prediction(Training) button to reopen the Data dialog. This time, select Expression (Simulation) and Expression (Training).

## Expression Simulation vs Expression Training results



This chart displays the predicted cost per room (MEDV/RM) for the Training partition as well as the synthetic data.

### ***LinReg\_Stored***

For information on Stored Model Sheets, in this example *LinReg\_Stored*, please refer to the “Scoring New Data” chapter that appears later in this guide.

# Scoring New Data

---

## Introduction

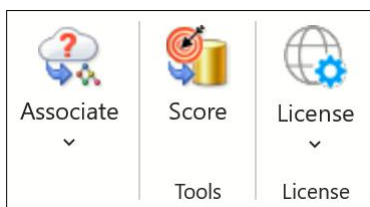
Analytic Solver Data Science provides a method for scoring new data in a database or worksheet with any of the prediction, classification, forecasting or transformation features.

---

## Scoring New Data

When Analytic Solver Data Science calculates prediction, classification, forecasting and transformation results, internal values and coefficients are generated and used in the computations. These values are saved to output sheets named, X\_Stored -- where X is the abbreviated name of the data science method. For example, the name given to the stored model sheet for Linear Regression is "LinReg\_Stored".

Note: In previous versions of XLMiner, this utility was a separate add-on application named XLMLCalc. Starting with XLMiner V12, this utility is included free of charge and can be accessed under Score in the Tools section of the XLMiner ribbon.



For example, assume the Linear Regression prediction method has just finished. The Stored Model Sheet (*LinReg\_Stored*) will contain the regression equation in PMML format. When the Score Test Data utility is invoked or when the *PsiPredict()* function is present (see below), Analytic Solver will apply this equation from the Stored Model Sheet to the test data.

Along with values required to generate the output, the Stored Model Sheet also contains information associated with the input variables that were present in the training data. The dataset on which the scoring will be performed should contain *at least* these original Input variables. Analytic Solver Data Science offers a “matching” utility that will match the Input variables in the training set to the variables in the new dataset so the variable names are not required to be identical in both data sets (training and test). See the sections below for more information on scoring to a database.

## Scoring New Data Example

This example illustrates how to score new data using a stored model sheet using output from a Multiple Linear Regression fitted model. In other words, the fitted Linear Regression model from the Predicting Housing Prices using MLR will be used to predict values for the MEDV variable for 10 *new* housing tracts. The new dataset may be found below. This procedure may be repeated using any stored model sheets generated using Analytic Solver Data Science.

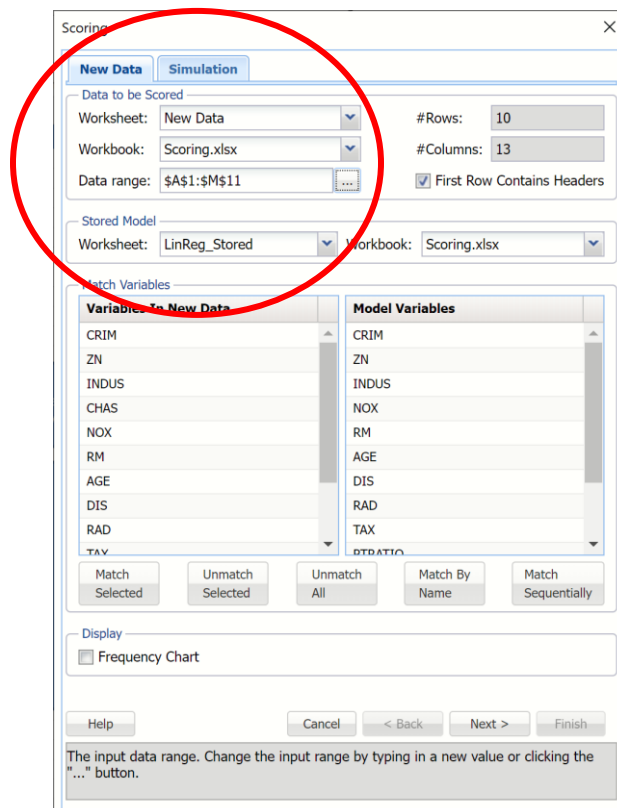
|    | A        | B  | C     | D    | E     | F     | G    | H       | I   | J   | K       | L      | M     |
|----|----------|----|-------|------|-------|-------|------|---------|-----|-----|---------|--------|-------|
| 1  | CRIM     | ZN | INDUS | CHAS | NOX   | RM    | AGE  | DIS     | RAD | TAX | PTRATIO | B      | LSTAT |
| 2  | 0.00487  | 16 | 2.25  | 0    | 0.65  | 5.454 | 60   | 3.75    | 1   | 265 | 15.3    | 378    | 4.5   |
| 3  | 1.2456   | 1  | 8     | 0    | 0.522 | 5.5   | 95   | 3.65    | 4   | 298 | 20      | 366.57 | 20.2  |
| 4  | 0.03495  | 77 | 3     | 0    | 0.824 | 8.65  | 18.5 | 4.511   | 2   | 265 | 17      | 410    | 2.05  |
| 5  | 0.15     | 26 | 4.98  | 0    | 0.354 | 3.25  | 60   | 5.7422  | 9   | 295 | 23.7    | 591.31 | 13.15 |
| 6  | 0.87038  | 1  | 13.28 | 0    | 0.734 | 7.854 | 66.3 | 6.5402  | 5   | 389 | 17.8    | 369.06 | 10.9  |
| 7  | 0.0866   | 1  | 8.29  | 0    | 0.544 | 6.741 | 56.2 | 9.543   | 1   | 267 | 19      | 399.6  | 9.61  |
| 8  | 10.587   | 0  | 11    | 0    | 0.745 | 5.872 | 71.3 | 7.2457  | 7   | 423 | 18.7    | 383.36 | 17.53 |
| 9  | 3.22158  | 1  | 18.95 | 0    | 0.56  | 9.436 | 74.9 | 1.78773 | 6   | 340 | 17.4    | 634.33 | 5.94  |
| 10 | 0.06426  | 0  | 5.4   | 0    | 0.15  | 0.866 | 44.7 | 3.5921  | 4   | 629 | 15.3    | 379.12 | 6.29  |
| 11 | 0.086703 | 84 | 2.51  | 0    | 0.44  | 2.742 | 8.33 | 9.73    | 3   | 239 | 16.2    | 329.2  | 6.26  |

Click **Help – Example Models** on the Data Science ribbon, then **Forecasting/Data Science Examples** and open the example file **Scoring.xlsx**.

*LinReg\_Stored* was generated while performing the steps in the “Multiple Linear Regression Prediction Method” chapter. See this chapter for details on performing a Multiple Linear Regression.

Scoring.xlsx also contains a *New Data* worksheet with 10 new records. Our goal is to score this new dataset to come up with a predicted housing price for each of the 10 new records.

Click **Score** on the Data Science ribbon. Under *Data to be scored*, confirm that *New Data* appears as the Worksheet, Scoring.xlsx as the Workbook, the Data range is **A1:M11** and *LinReg\_Stored* is selected in the *Worksheet* drop down menu under *Stored Model*.



Variables in the New Data may be matched with Variables in Stored Model using three easy techniques: by name, by sequence or manually.

If **Match By Name** is clicked, all similar named variables in the stored model sheet will be matched with similar named variables in the new dataset.

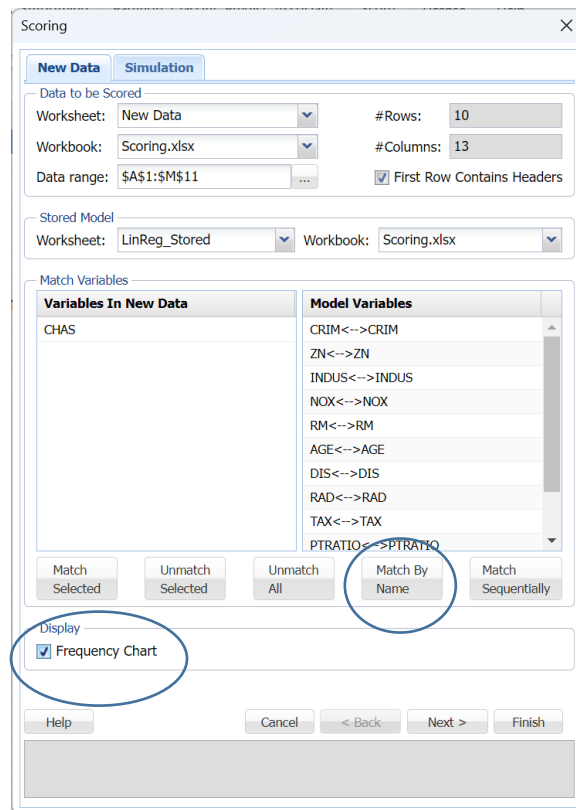
If **Match Sequentially** is clicked, the Variables in the stored model will be matched with the Variables in the new data in order that they appear in the two list boxes. For example, the variable *CRIM* from the new dataset will be matched with the variable *CRIM* from the stored model sheet, the variable *ZN* from the new data will be matched with the variable *ZN* from the stored model sheet and so on.

To manually map variables from the stored model sheet to the new data set, select a variable from the new data set in the *Variables in New Data* list box, then select the variable to be matched in the stored model sheet in the *Variables in Stored Model* list box, then click **Match**. For example to match the *CRIM* variable in the new dataset to the *CRIM* variable in the stored model sheet, select **CRIM** from the *Variables in New Data* list box, select **CRIM** from the stored model sheet in the *Variables in Stored Model* list box, then click **Match Selected** to match the two variables.

To unmatch all variables click **Unmatch all**. To unmatch two specific variables, select the matched variables, then click **Unmatch Selected**.

Click **Match By Name** to quickly match all variables in the fitted model with the variables in the new data.

Click the Frequency Chart checkbox, in the bottom left of the chart, to display frequency charts for all variables in the output.

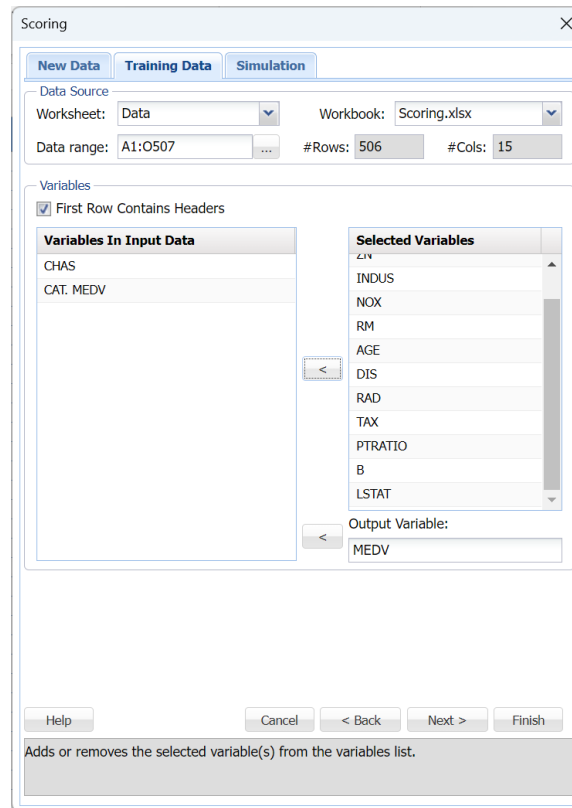


From here you can either click **Finish** to score the data on the New Data worksheet. Or, if you are scoring a classification or prediction model, you can click **Next** to advance to the Simulation tab where you can perform a complete risk analysis on the data.

Click **Next**.



Once the Simulate Response Prediction checkbox is selected, a new Training Data tab appears to the left of the Simulation tab. Use this tab to select the continuous variables to be included in the risk analysis. Make sure to select Worksheet: Data and Data range: A1:O507 within the Data Source section of the Training Data tab.



The same variables must be selected as what were selected previously when the stored model sheet was created.

Click the LinReg\_Output worksheet to see the variables selected for the original fitted model: CRIM, ZN, INDUS, NOX, RM, AGE, DIS, RAD, TAX, PTRATIO, B and LSTAT with MEDV as the output variable.

Click Next to return to the Simulation tab. Options for the risk analysis of the new data using the trained and validated machine learning model fit in the previous Predicting Housing Prices using MLR chapter, are set on the Simulation tab. This example uses the default of 100 simulated cases and enters an expression to calculate the median price per room in each census tract.

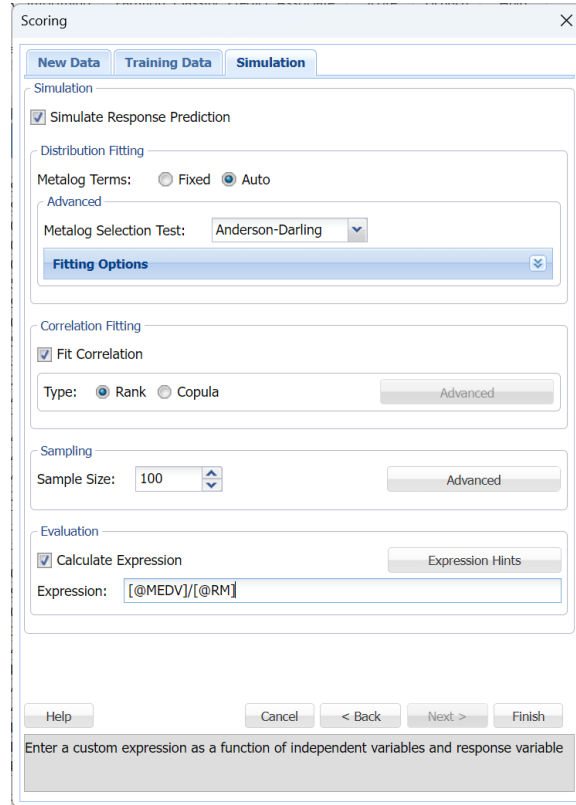
Enter the following formula for Expression:  $[@MEDV]/[@RM]$

A chart with the results of this expression as applied to both the new data and the training data will be inserted into the output.

All options on the Simulation tab are left at their default settings. Recall that Fitting Options (click the down arrow to open) controls the automated fitting of Metalog probability distributions to each feature in the dataset. By default, bounds in the dataset are used as bounds for the family of Metalog distributions. However, users can easily and quickly adjust or remove the lower and/or upper bounds. Correlation Fitting will, by default, construct a rank-order correlation matrix that includes all features, but choosing “Copula” activates the “Advanced” button, where types of copulas (Clayton, Frank, Gumbel, etc) may

be selected. For more information on each option, see the Generate Data chapter within the Data Science Reference Guide.

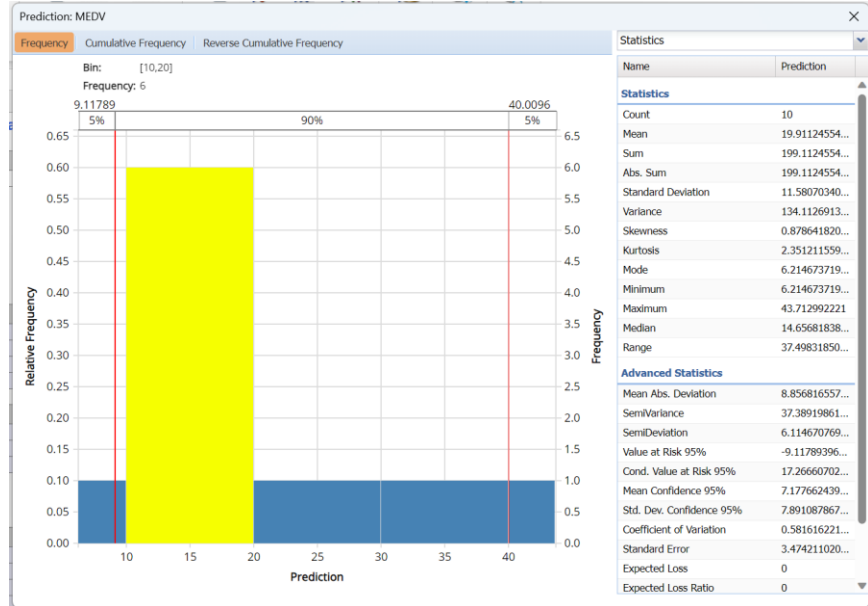
Click Finish to score the new data (from the New Data worksheet) and the synthetic data and to also perform a risk analysis on the original data (from the Data worksheet).



Two new worksheets, Scoring\_LinearRegression and Scoring\_Simulation, are inserted to the right of the LinReg\_Stored tab.

### ***Scoring\_LinearRegression***

The first thing you'll notice is the frequency chart that appears when the tab opens. This chart displays the frequency of the output variable, MEDV, for the new data dataset. Use the mouse to hover over any of the bars in the graph to populate the Bin and Frequency headings at the top of the chart.

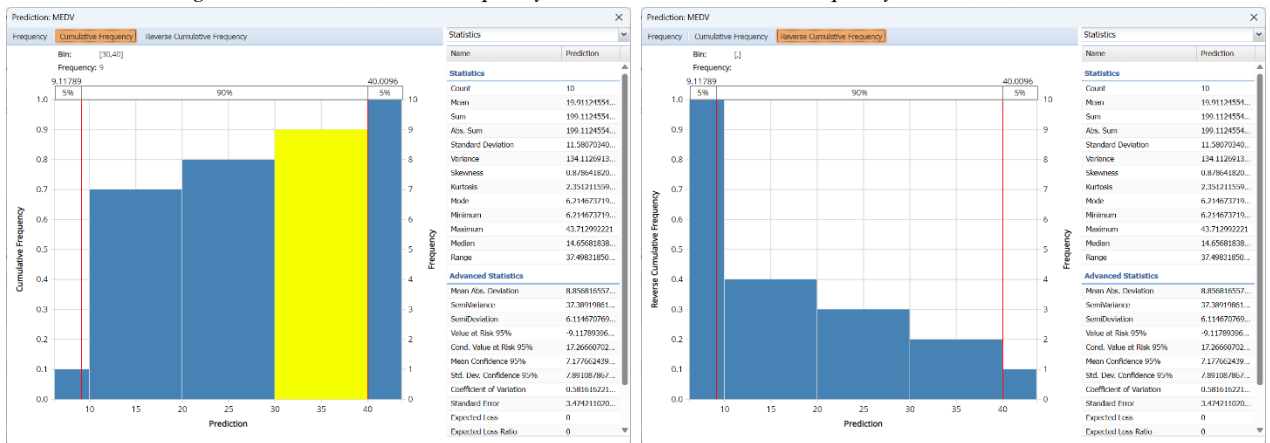


From this chart, you can see the number of records contained in the new data plus Advanced and Summary statistics.

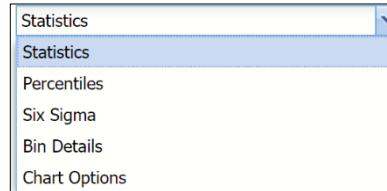
Red vertical lines appear at the 5% and 95% percentile values in all three charts (Frequency, Cumulative Frequency and Reverse Cumulative Frequency) effectively displaying the 90<sup>th</sup> confidence interval. The middle percentage is the percentage of all the variable values that lie within the 'included' area, i.e. the darker shaded area. The two percentages on each end are the percentage of all variable values that lie outside of the 'included' area or the "tails". Percentile values can be altered by moving either red vertical line to the left or right.

Click Cumulative Frequency and Reverse Cumulative Frequency tabs to see the Cumulative Frequency and Reverse Cumulative Frequency charts, respectively.

*Prediction dialog shown with Cumulative Frequency and Reverse Cumulative Frequency charts*



Click the down arrow next to Statistics to view Percentiles, Six Sigma metrics, details related to the histogram bins in the chart and also to change the look of the chart using Chart Options.



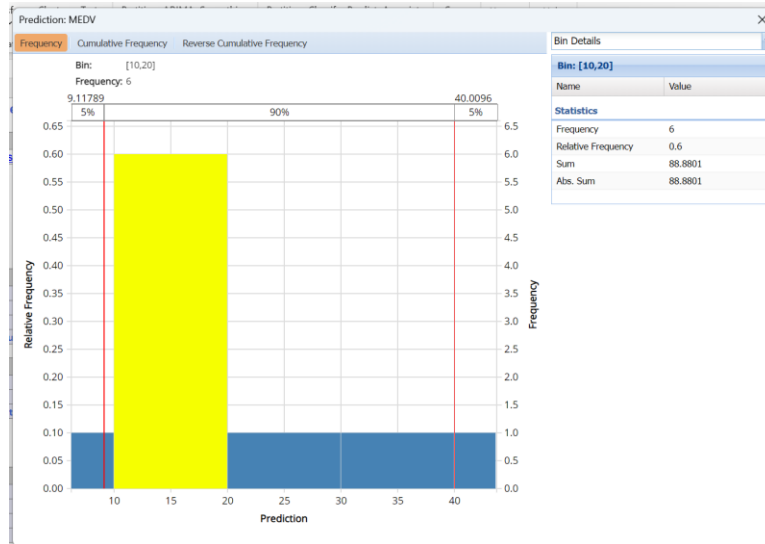
*Percentiles and Six Sigma views*

| Percentiles        |                | Six Sigma        |                 |
|--------------------|----------------|------------------|-----------------|
| Name               | Prediction     | Name             | Prediction      |
| <b>Percentiles</b> |                | <b>Six Sigma</b> |                 |
| 1%                 | 6.795317769... | Cp               | 0.444586162...  |
| 2%                 | 7.375961818... | Cpk              | 0.310670581...  |
| 3%                 | 7.956605867... | Cpk, lower       | 0.310670581...  |
| 4%                 | 8.537249916... | Cpk, upper       | 0.578501744...  |
| 5%                 | 9.117893966... | Cpm              | 0.223521537...  |
| 6%                 | 9.698538015... | PPM              | 216990.9057...  |
| 7%                 | 10.27918206... | PPM, lower       | 175665.2279...  |
| 8%                 | 10.85982611... | PPM, upper       | 41325.67776...  |
| 9%                 | 11.44047016... | K                | 0.301214010...  |
| 10%                | 12.02111421... | Lower Bound      | -49.57297486... |
| 11%                | 12.60175826... | ProbDefectShift  | 0.216990905...  |
| 12%                | 12.73140071... | ProbLowerShift   | 0.175665227...  |
| 13%                | 12.80466797... | ProbUpperShift   | 0.041325677...  |
| 14%                | 12.87793522... | Sigma Level      | 0.782396123...  |
| 15%                | 12.95120248... | Upper Bound      | 89.39546596...  |
| 16%                | 13.02446973... | Yield            | 0.783009094...  |
| 17%                | 13.09773699... | Z Lower          | 0.932011744...  |
| 18%                | 13.17100424... | Z Min            | 0.932011744...  |
| 19%                | 13.24427150... | Z Upper          | 1.735505232...  |
| 20%                | 13.31753875... |                  |                 |
| 21%                | 13.39080601... |                  |                 |
| 22%                | 13.46407326... |                  |                 |
| 23%                | 13.50409678... |                  |                 |
| 24%                | 13.53462208... |                  |                 |
| 25%                | 13.56514738... |                  |                 |

Select Bin Details to view details pertaining to each bin in the chart.

Hover over the bars in the chart to display important bin statistics such as frequency, relative frequency, sum and absolute sum.

*Bin Details view*



- Frequency is the number of observations assigned to the bin.
- Relative Frequency is the number of observations assigned to the bin divided by the total number of observations.
- Sum is the sum of all observations assigned to the bin.
- Absolute Sum is the sum of the absolute value of all observations assigned to the bin, i.e. |observation 1| + |observation 2| + |observation 3| + ...

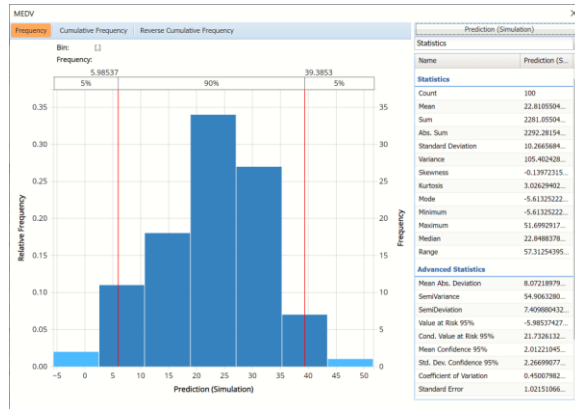
As discussed above, see the Generate Data section of the Exploring Data chapter (within the Analytic Solver Reference Guide) for an in-depth discussion of this chart as well as descriptions of all statistics, percentiles, bin details and six sigma indices.

Click the X in the upper right hand corner to close the frequency chart. You can reopen the chart simply by clicking on a different worksheet tab and then clicking back to the Scoring\_LinearRegression tab.

The results from scoring can be found under: Scoring. These are the mean predicted prices for the MEDV variable for each of the ten new records.

|    | B              | C                | D                       |
|----|----------------|------------------|-------------------------|
| 54 | <b>Scoring</b> |                  |                         |
| 55 |                |                  |                         |
| 56 |                | <b>Record ID</b> | <b>Prediction: MEDV</b> |
| 57 |                | Record 1         | 24.82155793             |
| 58 |                | Record 2         | 14.54256942             |
| 59 |                | Record 3         | 35.48326928             |
| 60 |                | Record 4         | 13.8195249              |
| 61 |                | Record 5         | 19.6001715              |
| 62 |                | Record 6         | 14.77106735             |
| 63 |                | Record 7         | 6.21467372              |
| 64 |                | Record 8         | 43.71299222             |
| 65 |                | Record 9         | 12.66627427             |
| 66 |                | Record 10        | 13.48035488             |

Click the 2<sup>nd</sup> output tab, Scoring\_Simulation. You'll notice a similar chart appear. This chart displays a frequency histogram of the predicted values for the synthetic data.



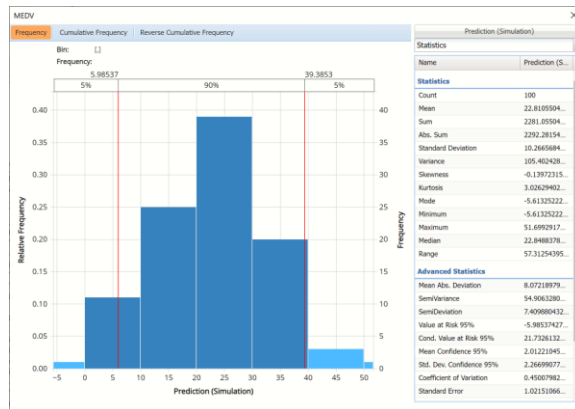
## Scoring\_Simulation

As discussed in the earlier chapters Automated Risk Analysis of ML Models and Fitting the Best Model and Predicting Housing Prices using MLR, Analytic Solver Data Science generates a new output worksheet, Scoring\_Simulation, when *Simulate Response Prediction* is selected on the Simulation tab of the Scoring dialog.

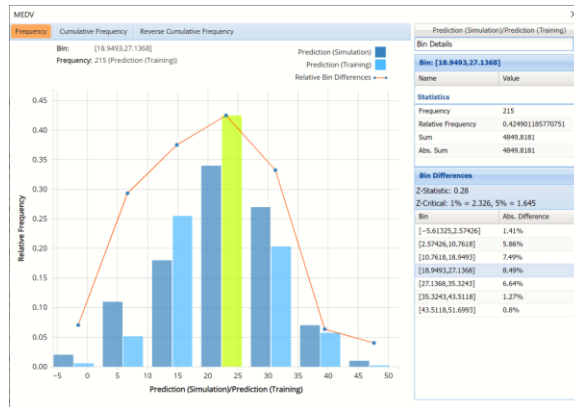
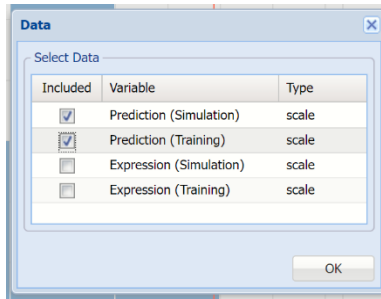
From this chart, users can view charts containing histograms of the synthetic data, the training data and the results of the Expression applied to both datasets, if populated in the dialog.

This chart is similar to the chart on the Scoring\_LinearRegression output sheet in that it also contains the same three charts described above (Frequency, Cumulative Frequency and Reverse Cumulative Frequency) and the same views: Statistics, Percentiles, Six Sigma and Bin Details. Click each chart tab to view the selected chart and click the down arrow next to Statistics to change the view.

By default, this chart first opens to the Prediction (Simulation) view, which displays a histogram of the predicted values in the simulation (or synthetic) data.



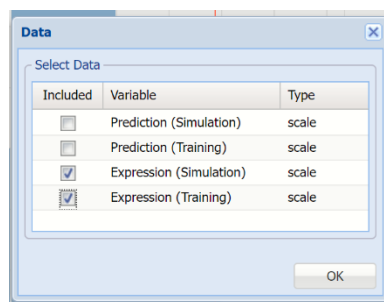
Click Prediction (Simulation) to open the Data dialog. Select Prediction (Simulation) and Prediction (Training) to view the predicted values in both the training and synthetic datasets.



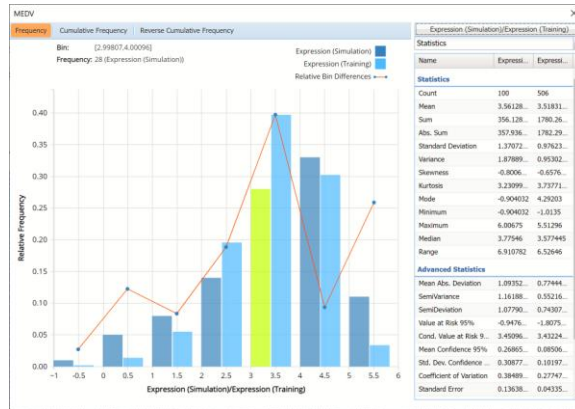
The dark blue bars display the frequencies of the MEDV variable in the synthetic data, or Prediction (Simulation), and the lighter blue bars display the frequencies of the MEDV variable in the original data, or Prediction (Training). The synthetic data is predicting a slightly higher number of homes in the \$27,000 to \$35,000 (remember these are 1940's housing prices!) range and fewer number of homes in the \$19,000 to \$27,000 range.

Notice that red curve which connects the relative Bin Differences for each bin. As discussed in the previous chapters, Bin Differences are computed based on the frequencies of records which predictions fall into each bin. For example, consider the highlighted bin in the screenshot above  $[x_0, x_1] = [18.949, 27.137]$ . This bin contains 34 records in the synthetic data and 215 records in the training data. The relative frequency of the Simulation data is  $34/100 = 34\%$  and the relative frequency of the Training data is  $215/506 = 42.5\%$ . Hence the Absolute Difference (in frequencies) is  $42.5 - 34 = 8.5\%$ .

Click back to the Data dialog and select Expression (Simulation) and Expression (Training) to view the results of the expression for both datasets.



Recall the expression:  $[@MEDV]/[@RM]$  which is simply the median value of houses in each housing tract divided by the average number of rooms in each dwelling, or a rudimentary calculation of the price per room.



The highlighted bar above shows 28 records in the synthetic data are contained in the Bin: [2.998, 4.00] meaning that there are 28 records in the synthetic data where the value of the expression falls within this interval.

Click the X in the upper right hand corner to close the dialog to view the simulated data for each variable. (To reopen, simply click another tab and then click back to the Scoring\_Simulation tab. )

| Record    | Expression  | MEDV     | CRIM     | ZN       | INDUS    | NOX      | RM       | AGE      | DIS         | RAD      | TAX      | PTRATIO    | B        | LSTAT     |
|-----------|-------------|----------|----------|----------|----------|----------|----------|----------|-------------|----------|----------|------------|----------|-----------|
| Record 1  | 0.819310585 | 3.660767 | 1.452193 | 2.13E-05 | 21.78657 | 0.527574 | 4.468106 | 59.24304 | 4.941903636 | 2.85584  | 528.6146 | 21.4121003 | 340.5634 | 20.68961  |
| Record 2  | 1.528528516 | 8.402374 | 1.1002   | 5.84E-06 | 7.934002 | 0.553812 | 5.497034 | 86.12905 | 3.843026692 | 9.132196 | 461.2002 | 21.8983229 | 396.8668 | 25.37774  |
| Record 3  | 5.41600715  | 42.33422 | 0.074801 | 99.64343 | 1.267204 | 0.42239  | 7.8165   | 13.39563 | 7.874093364 | 23.92055 | 342.8905 | 16.9064806 | 396.3879 | 3.3400577 |
| Record 4  | 3.709917413 | 20.60682 | 0.252207 | 4.76E-05 | 18.82521 | 0.53457  | 5.554523 | 41.79345 | 4.651889949 | 23.99917 | 690.9174 | 18.7314986 | 387.9219 | 11.42005  |
| Record 5  | 4.09593977  | 33.64456 | 0.119144 | 7.43E-11 | 5.887465 | 0.536573 | 8.214125 | 99.99866 | 3.542683057 | 1.001093 | 238.182  | 19.100286  | 396.8938 | 3.572565  |
| Record 6  | 2.006136786 | 14.78957 | 7.050132 | 3.08E-09 | 22.3977  | 0.852757 | 7.372163 | 99.99996 | 1.290690971 | 23.99955 | 619.8849 | 20.6707419 | 396.5248 | 30.50952  |
| Record 7  | 5.764995761 | 37.8151  | 0.159871 | 3.58E-05 | 7.077398 | 0.474599 | 6.559432 | 39.89342 | 3.470775246 | 23.88007 | 268.2903 | 17.3244989 | 396.9    | 5.395022  |
| Record 8  | 3.459166788 | 22.39342 | 0.013615 | 0.001487 | 4.003189 | 0.391686 | 6.473647 | 8.935999 | 10.71283449 | 1.006276 | 213.0946 | 14.1146607 | 396.459  | 3.731889  |
| Record 9  | 3.327596575 | 20.87533 | 0.015485 | 0.000548 | 3.797742 | 0.508287 | 6.273395 | 44.91892 | 7.505854413 | 10.80255 | 245.6237 | 17.6875647 | 396.9    | 11.48939  |
| Record 10 | 0.780749647 | 3.149372 | 1.595809 | 1.2E-09  | 26.76259 | 0.53751  | 4.03378  | 89.46242 | 4.144841901 | 4.491779 | 444.2761 | 21.2071115 | 358.2376 | 26.53406  |
| Record 11 | 4.751383469 | 32.56936 | 0.480823 | 1.43E-06 | 6.711089 | 0.537735 | 6.85471  | 56.81309 | 2.835185846 | 23.99939 | 481.2724 | 13.7266116 | 396.6464 | 15.12485  |
| Record 12 | 3.854146602 | 25.48643 | 0.108477 | 2.33E-05 | 5.609579 | 0.536499 | 6.612731 | 51.03386 | 4.799208932 | 23.98856 | 571.5908 | 14.8041738 | 0.408557 | 8.044393  |

## Using Data Science Psi Functions in Excel

Analytic Solver Data Science utilizes XML PMML format to store the supported models and use them for “scoring” (classifying, predicting, forecasting, transforming) new data using four new generic scoring functions: **PsiPredict()**, **PsiForecast()**, **PsiTransform()** and **PsiPosteriors()**. PsiPredict() and PsiForecast() provide previously available functionality plus new additional functionality such as storing and scoring ensemble models with any available weak learner and computing the fitted values for new time series data. PsiTransform() and PsiPosteriors() provide new functionality and the availability of new models for storing or scoring.

| PSI Scoring Function | Description                                                                                                                                                 |
|----------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------|
| PsiPredict()         | Predicts the target for input data using a Classification or Regression model and computes the fitted values for a Time Series model stored in PMML format. |
| PsiForecast()        | Computes the forecasts for the input data using a Time Series model stored in PMML format.                                                                  |
| PsiPosteriors()      | Computes the posterior probabilities for the input data                                                                                                     |



|                |                                                                               |
|----------------|-------------------------------------------------------------------------------|
|                | using a Classification model stored in PMML format.                           |
| PsiTransform() | Transforms the input data using a Transformation model stored in PMML format. |

It's possible to score data with a prediction or classification method or perform a time series forecast manually (without the need to click the Score icon on the ribbon) by entering a Psi Solver function into an Excel cell as an array.

## Scoring Data Using Psi Functions

Using the example above, click over to the *PsiPredict* worksheet and select a blank cell on the worksheet. Enter "`=PsiPredict(LinReg_Stored!B12:B69,'New Data'!A1:L11)`". The formula result will "spill" into the cells below.

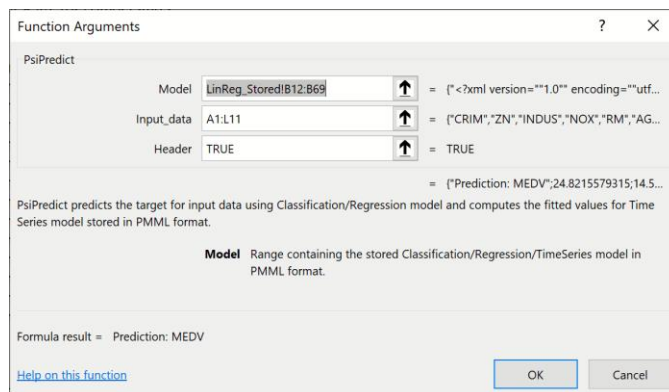
Note: If using a version of Excel that does not support Dynamic Arrays, select cells N1:N11, enter the formula and press SHIFT + CTRL + ENTER to enter the formula as an array into all 10 cells (N2:N11).

The first argument, `LinReg_Stored!B12:B69`, is the range of cells used by Analytic Solver Data Science to store the linear regression model on the *LinReg\_Stored* worksheet. Clearly this data range will change as the classification or prediction method changes and as the number of features included in the dataset changes.

The second argument, `New Data!A1:L11`, is the range containing the new data on the *New Data* worksheet. The new data must contain at least one row of data containing the same number of features (or columns) as the data used to create the model. In this example, we included 12 features in the linear regression model: CRIM, ZN, INDUS, NOX, RM, AGE, DIS, RAD, TAX, PTRATIO, B & LSTAT. As a result, our new data also contains these same 12 features. We could have performed this prediction on only one row of new data, `NewData!A2:L2`, but choose to use all 10 rows.

Note: If there are new or missing features in the new dataset, the Psi Data Science function will return #VALUE.

It's also possible to enter this formula using the Insert Function dialog by clicking Formulas – Insert Function, select PSI Data Science for Category, then PsiPredictMLR.



Note: The Insert Function dialog is not supported when using Data Science Cloud. To use this function or any other Psi function, simply type the function directly into the cell.

The scoring results are shown in the screenshot below in the N column.

|    | A        | B  | C     | D     | E     | F    | G       | H   | I   | J       | K      | L     | M | N                | O |
|----|----------|----|-------|-------|-------|------|---------|-----|-----|---------|--------|-------|---|------------------|---|
| 1  | CRIM     | ZN | INDUS | NOX   | RM    | AGE  | DIS     | RAD | TAX | PTRATIO | B      | LSTAT |   | Prediction: MEDV |   |
| 2  | 0.00487  | 16 | 2.25  | 0.85  | 5.454 | 60   | 3.75    | 1   | 285 | 15.3    | 378    | 4.5   |   | 24.8215579       |   |
| 3  | 1.2456   | 1  | 8     | 0.522 | 5.5   | 95   | 3.65    | 4   | 298 | 20      | 366.57 | 20.2  |   | 14.5425694       |   |
| 4  | 0.03495  | 77 | 3     | 0.824 | 8.65  | 18.5 | 4.511   | 2   | 265 | 17      | 410    | 2.05  |   | 35.4832693       |   |
| 5  | 0.15     | 26 | 4.98  | 0.354 | 3.25  | 60   | 5.7422  | 9   | 295 | 23.7    | 591.31 | 13.15 |   | 13.8195249       |   |
| 6  | 0.87038  | 1  | 13.28 | 0.734 | 7.854 | 66.3 | 6.5402  | 5   | 389 | 17.8    | 369.06 | 10.9  |   | 19.6001715       |   |
| 7  | 0.0866   | 1  | 8.29  | 0.544 | 6.741 | 56.2 | 9.543   | 1   | 267 | 19      | 399.6  | 9.61  |   | 14.7710674       |   |
| 8  | 10.587   | 0  | 11    | 0.745 | 5.672 | 71.3 | 7.2457  | 7   | 423 | 18.7    | 383.36 | 17.53 |   | 6.21467372       |   |
| 9  | 3.22158  | 1  | 18.95 | 0.56  | 9.436 | 74.9 | 1.78773 | 6   | 340 | 17.4    | 634.33 | 5.94  |   | 43.7129922       |   |
| 10 | 0.06426  | 0  | 5.4   | 0.15  | 0.866 | 44.7 | 3.5921  | 4   | 629 | 15.3    | 379.12 | 8.29  |   | 12.6662743       |   |
| 11 | 0.086703 | 84 | 2.51  | 0.44  | 2.742 | 8.33 | 9.73    | 3   | 239 | 16.2    | 329.2  | 6.26  |   | 13.4803549       |   |
| 12 |          |    |       |       |       |      |         |     |     |         |        |       |   |                  |   |

The PsiPredict() function is interactive meaning that if a variable value is changed, for example the first LSTAT value in cell L2 changes from 4.5 to 9.5, the Predicted Value in cell N2 will immediately update to reflect a new predicted value.

The remaining PSI Data Science functions can be used in the same way using models from their respective stored model sheets. See below for specifications for PsiPredict(), PsiPosteriors() and PsiTransform(). See the section below for information on PsiForecast().

## PsiPredict()

`PsiPredict(Model, Input_Data, [Header])`

Predicts the response, target, output or dependent variable for `Input_Data` whether it is continuous (Regression) or categorical (Classification) when the model is stored in PMML format. In addition, this function also computes the fitted values for a Time Series model when the model is stored in PMML format.

**Model:** Range containing the stored Classification, Regression or TimeSeries model in PMML format.

**Input\_Data:** Range containing the new data for computing predictions.

Range must contain a header row with column names and at least one row of data containing the exact same features (or columns) as the data used to create the model.

**Header:** If True, a heading will be inserted above the forecasted values. If omitted or false, a heading will not appear.

The contents of the Dynamic Array will "spill" down the column. If a nonblank cell is "blocking" the contents of the Dynamic Array, PsiForecast() will return #SPILL until such time as the blockage is removed. Use the optional `numForecasts` argument to specify the number of forecasts in the Dynamic Array. If not present, one forecast will be returned.

**Output:** A single column containing the header (if the header argument is set to TRUE) and predicted/fitted values for each record in `Input_Data`.

To know if the result of the prediction is continuous or categorical, you must know what kind of model you are passing as an argument to the scoring function – if you previously fitted the classification model and are now predicting the new feature vectors, you should expect to get the compatible categorical response. On the other hand, you should expect the continuous response from the new data prediction when using a fitted regression model. In previous versions, the user was expected to know the exact type model, such as Multiple Linear Regression or Discriminant Analysis, to know what kind of output will be produced, whereas in V2017 and later, it is sufficient to know whether you're pointing to a classification or regression model in order to determine the type of

the response. Note: If the user intends to use an “unknown” model for scoring, the stored worksheets contain the complete information about the model including several clear indications of the model type and data dictionaries with the types of features and response.

In addition, PsiPredict() can compute the *fitted values* for the new time series based on the provided Time Series model. Unlike future-looking forecasting, provided by PsiForecast(), PsiPredict() computes a model *prediction* for each observation in the provided new time series.

*Supported Models:*

- Classification:
  - Linear Discriminant Analysis
  - Logistic Regression
  - K-Nearest Neighbors
  - Classification Tree
  - Naïve Bayes
  - Neural Network
  - Random Trees
  - Bagging (with any supported weak learner)
  - Boosting (with any supported weak learner)
- Regression:
  - Logistic Regression
  - K-Nearest Neighbors
  - Neural Network
  - Bagging (with any supported weak learner)
  - Boosting (with any supported weak learner)
- Time Series (fitted values)
  - ARIMA
  - Exponential Smoothing
  - Double Exponential Smoothing
  - Holt-Winters Smoothing

*Previous related Psi Scoring functions:*

- Classification: PsiClassifyLR, PsiClassifyDA, PsiClassifyCT, PsiClassifyNB, PsiClassifyNN, PsiClassifyCTEnsemble, PsiClassifyNNEnsemble
- Regression: PsiPredictMLR, PsiPredictRT, PsiPredictNN, PsiPredictNNEnsemble, PsiPredictRTEnsemble

| <b>Prediction/Classification/Time Series Algorithm</b> | <b>Stored Model Sheet</b>                                |
|--------------------------------------------------------|----------------------------------------------------------|
| Linear Discriminant Analysis Classification            | DA_Stored                                                |
| Logistic Regression Classification                     | LogReg_Stored                                            |
| k-Nearest Neighbors Classification                     | KNNC_Stored                                              |
| Classification Trees                                   | CT_Stored                                                |
| Naïve Bayes Classification                             | NB_Stored                                                |
| Neural Networks Classification                         | NNC_Stored                                               |
| Ensemble Methods for Classification                    | CBoosting_Stored<br>CBagging_Stored<br>CRandTrees_Stored |

|                                 |                                                                              |
|---------------------------------|------------------------------------------------------------------------------|
| Linear Regression               | LinReq_Stored                                                                |
| k-Nearest Neighbors Regression  | KNNP_Stored                                                                  |
| Regression Tree                 | RT_Stored                                                                    |
| Neural Network Regression       | NNP_Stored                                                                   |
| Ensemble Methods for Regression | RBoosting_Stored<br>RBagging_Stored<br>RRandTrees_Stored                     |
| ARIMA                           | ARIMA_Stored                                                                 |
| Exponential Smoothing           | Expo_Stored                                                                  |
| Double Exponential Smoothing    | DoubleExpo_Stored                                                            |
| Moving Average Smoothing        | MovingAvg_Stored                                                             |
| Holt Winters Smoothing          | MultHoltWinters_Stored<br>AddHoltWinters_Stored<br>NoTrendHoltWinters_Stored |

## PsiPosteriors()

---

`PsiPosteriors(Model, Input_Data, [Header])`

Computes the posterior probabilities for `Input_Data` using a Classification model stored in PMML format.

**Model:** Range containing the stored Classification model in PMML format.

**Input\_Data:** Range containing the new data for computing posterior probabilities. Range must contain a header with column names and at least one row of data containing the exact same features (or columns) as the data used to create the model.

**Header:** If True, a heading is inserted in the output above the forecasted values. If False or omitted, a heading is not inserted into the output.

In Data Science Cloud and in new versions of desktop Excel, `PsiPosterior()` returns a [Dynamic Array](#) (see Note in section heading, above) To use this function in the Cloud, you need only enter the Psi function in one cell as a normal function, i.e., not as a control array. The contents of the Dynamic Array will "spill" down the column. If a nonblank cell is "blocking" the contents of the Dynamic Array, `PsiForecast()` will return #SPILL until such time as the blockage is removed. Use the optional `numForecasts` argument to specify the number of forecasts in the Dynamic Array. If not present, one forecast will be returned.

**Output:** Multiple columns containing a header with class labels and estimated posterior probabilities for each class label for all records in `Input_Data`.

*Supported Models:*

- Classification:
  - Linear Discriminant Analysis
  - Logistic Regression

- K-Nearest Neighbors
- Classification Tree
- Naïve Bayes
- Neural Network
- Random Trees
- Bagging (with any supported weak learner)
- Boosting (with any supported weak learner)

*Previous related Psi Scoring functions: N/A*

| <b>Classification Algorithm</b>             | <b>Stored Model Sheet</b>                                |
|---------------------------------------------|----------------------------------------------------------|
| Linear Discriminant Analysis Classification | DA_Stored                                                |
| Logistic Regression Classification          | LogReg_Stored                                            |
| k-Nearest Neighbors Classification          | KNNC_Stored                                              |
| Classification Trees                        | CT_Stored                                                |
| Naïve Bayes Classification                  | NB_Stored                                                |
| Neural Networks Classification              | NNC_Stored                                               |
| Ensemble Methods for Classification         | CBoosting_Stored<br>CBagging_Stored<br>CRandTrees_Stored |

## **PsiTransform()**

`PsiTransform(Model, Input_Data, [Header])`

Transforms the `Input_Data` using a Transformation model stored in PMML format.

**Model:** Range containing the stored Transformation model in PMML format.

**Input\_Data:** Range containing the new data for transformation. Range must contain a header with column names and at least one row of data containing the exact same features (or columns) as the data used to create the model.

**Header:** If True, a heading is inserted above the forecasted values. If False or omitted, a heading is not inserted.

In Data Science Cloud and in newer versions of desktop Excel, `PsiTransform()` returns a Dynamic Array (see Note in section heading, above). To use this function in the Cloud, you need only enter the Psi function in one cell as a normal function, i.e., not as a control array. The contents of the Dynamic Array will "spill" down the column. If a nonblank cell is "blocking" the contents of the Dynamic Array, `PsiForecast()` will return #SPILL until such time as the blockage is removed. Use the optional *numForecasts* argument to specify the number of forecasts in the Dynamic Array. If not present, one forecast will be returned.

**Output:** One or multiple columns containing a header and transformed data.

*Supported Models:*

- Transformation:
  - Rescaling
- Text Science

- TF-IDF Vectorization (input data – text variable with the corpus of documents)
- LSA Concept Extraction (input data – term-document matrix, where columns represent terms and rows represent documents)

*Previous related Psi Scoring functions: N/A*

| Algorithm    | Stored Model Sheet         |
|--------------|----------------------------|
| Rescaling    | Rescaling_Stored           |
| Text Science | TFIDF_Stored<br>LSA_Stored |

## Time Series Forecasting

Starting with version 2014-R2, Analytic Solver Data Science includes the ability to forecast a future point in a time series in one of your spreadsheet formulas (without using the Score button on the Ribbon) using a PsiForecast() function in conjunction with a model created using ARIMA or one of our smoothing methods (Exponential, Double Exponential, Moving Average, or Holt Winters).

PsiForecast() is similar to the previous PSIForecastXXX functions supported in V2014, 2015, and 2016: it will compute future-looking forecasts based on the fitted model, using the provided new time series observations as initial points. The result of PsiForecast() can be deterministic, if the Simulate argument is FALSE, or non-deterministic, if the Simulate argument is TRUE– in which case the forecasts are adjusted with random normally distributed errors, defined by the forecasts’ statistics.

Open the Airpass.xlsx example dataset by clicking Help – Examples on the Data Science ribbon, then clicking Forecasting/Data Science Examples. This example dataset includes International Airline Passenger Information by month for years 1949 – 1960. Since the number of airline passengers increases during certain times of the year, for example Spring, Summer, and in the month of December, we can say that this dataset includes “seasonality”.

First, we will partition this dataset into two datasets: a training dataset and a validation dataset. We’ll use the training dataset to create the ARIMA model and then we’ll apply the model to the validation dataset to forecast six future data points, or one half year of data.

Click Partition in the Time Series section of the Data Science ribbon to open the Time Series Partition Data dialog. Select Passengers for the Variables in the Partition Data and Month for the Time Variable.

## Time Series Forecasting

Starting with version 2014-R2, Analytic Solver Data Science includes the ability to forecast a future point in a time series in one of your spreadsheet formulas (without using the Score button on the Ribbon) using a PsiForecast() function in conjunction with a model created using ARIMA or one of our smoothing methods (Exponential, Double Exponential, Moving Average, or Holt Winters).

PsiForecast() is similar to the previous PSIForecastXXX functions supported in V2014, 2015, and 2016: it will compute future-looking forecasts based on the fitted model, using the provided new time series observations as initial points.

The result of PsiForecast() can be deterministic, if the Simulate argument is FALSE, or non-deterministic, if the Simulate argument is TRUE— in which case the forecasts are adjusted with random normally distributed errors, defined by the forecasts’ statistics.

Open the Airpass.xlsx example dataset by clicking Help – Examples on the Data Science ribbon, then clicking Forecasting/Data Science Examples. This example dataset includes International Airline Passenger Information by month for years 1949 – 1960. Since the number of airline passengers increases during certain times of the year, for example Spring, Summer, and in the month of December, we can say that this dataset includes “seasonality”.

First, we will partition this dataset into two datasets: a training dataset and a validation dataset. We’ll use the training dataset to create the ARIMA model and then we’ll apply the model to the validation dataset to forecast six future data points, or one half year of data.

Click Partition in the Time Series section of the Data Science ribbon to open the Time Series Partition Data dialog. Select Passengers for the Variables in the Partition Data and Month for the Time Variable.

## Time Series Forecasting

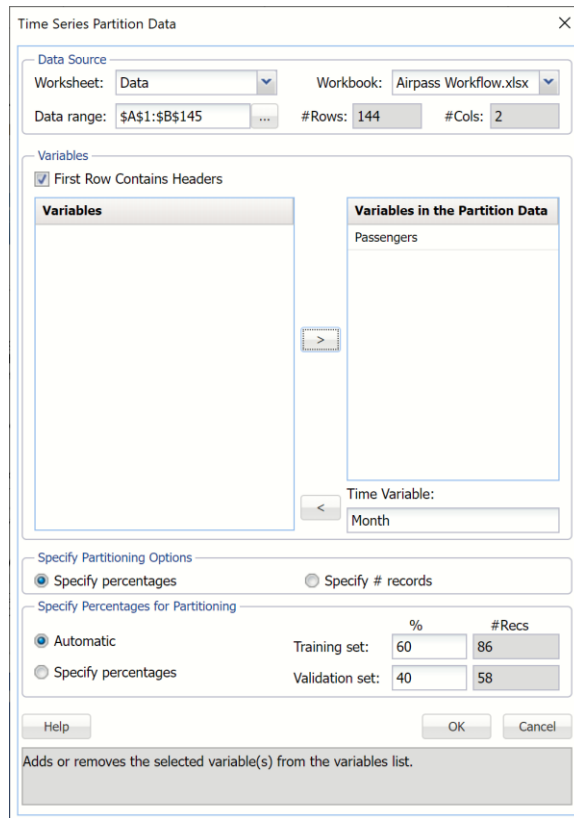
Starting with version 2014-R2, Analytic Solver Data Science includes the ability to forecast a future point in a time series in one of your spreadsheet formulas (without using the Score button on the Ribbon) using a PsiForecast() function in conjunction with a model created using ARIMA or one of our smoothing methods (Exponential, Double Exponential, Moving Average, or Holt Winters).

PsiForecast() is similar to the previous PSIForecastXXX functions supported in V2014, 2015, and 2016: it will compute future-looking forecasts based on the fitted model, using the provided new time series observations as initial points. The result of PsiForecast() can be deterministic, if the Simulate argument is FALSE, or non-deterministic, if the Simulate argument is TRUE— in which case the forecasts are adjusted with random normally distributed errors, defined by the forecasts’ statistics.

Open the Airpass.xlsx example dataset by clicking Help – Examples on the Data Science ribbon, then clicking Forecasting/Data Science Examples. This example dataset includes International Airline Passenger Information by month for years 1949 – 1960. Since the number of airline passengers increases during certain times of the year, for example Spring, Summer, and in the month of December, we can say that this dataset includes “seasonality”.

First, we will partition this dataset into two datasets: a training dataset and a validation dataset. We’ll use the training dataset to create the ARIMA model and then we’ll apply the model to the validation dataset to forecast six future data points, or one half year of data.

Click Partition in the Time Series section of the Data Science ribbon to open the Time Series Partition Data dialog. Select Passengers for the Variables in the Partition Data and Month for the Time Variable.



Click **OK** to accept the defaults for Specify Partitioning Options and Specify Percentages for Partitioning. Recall that when a time series dataset is partitioned, the dataset is partitioned sequentially. Therefore, 60% or the first 86 records, will be assigned to the training dataset and the remaining 40%, or 58 records, will be assigned to the validation dataset. (For more information on partitioning a time series dataset, see the previous chapter Exploring a Time Series Dataset.)

The TSPartition worksheet will be inserted into the Model tab of the Analytic Solver task pane under Transformations – Time Series Partition. Recall the steps needed to produce the forecast. Click ARIMA -- ARIMA to open the ARIMA dialog. Month has been pre selected as the Time variable. Select Passengers as the Selected variable.

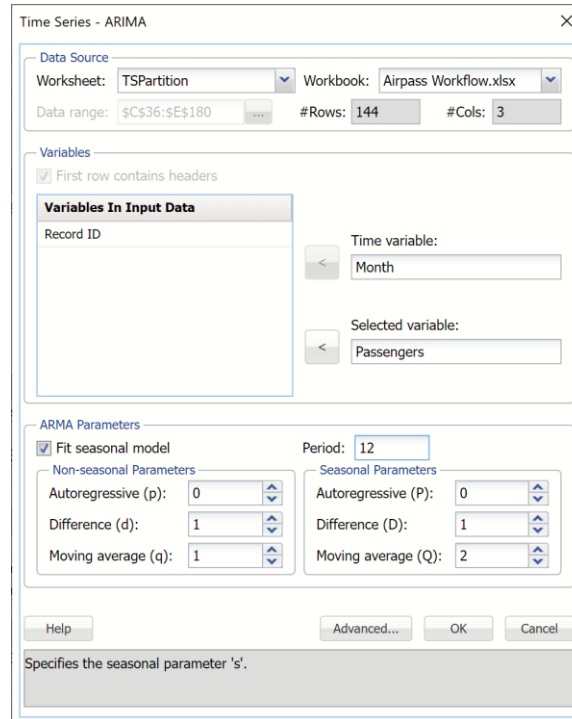
This example will use a SARIMA model, or Seasonal Autoregressive Integrated Moving Average model, to predict the next six datapoints in the dataset. (For more information on this type of time series model, please see the earlier chapter, “Exploring a Time Series Dataset.”) A seasonal ARIMA model requires 7 parameters, 3 nonseasonal (autoregressive (p), integrated (d), and moving average (q)), 3 seasonal (autoregressive (P), integrated (D), and moving average (Q)), and period. Each parameter must be a non-negative integer.

Selecting appropriate values for p, d, q, P, D, Q and period is beyond the scope of this User Guide. Consequently, this example will use a well documented SARIMA model with parameters  $p = 0$ ,  $d = 1$ ,  $q = 1$ ,  $P = 0$ ,  $D = 1$ ,  $Q = 2$  and period (P) = 12. Please refer to the classic time series analysis text **Time Series Analysis: Forecasting and Control** written by George Box and Gwilym Jenkins for more information on parameter selection.

Select Fit seasonal model and enter 12 for Period since it takes a full 12 months for the seasonal pattern to repeat. Set the Non-seasonal Parameters as



Autoregressive (p) = 0, Difference (d) = 1, Moving Average (q) = 1 and the Seasonal Parameters as Autoregressive (P) = 0, Difference (D) = 1, and Moving Average (Q) = 2.



Click **OK** to create the SARIMA model.

ARIMA\_Output will be inserted into the Model tab of the task pane under Reports – ARIMA. This output contains the Training Error Measures and Fitted Model Statistics. (For more information on this report, please see the chapter Exploring a Time Series Dataset within the Analytic Solver Data Science Reference Guide.) ARIMA\_Stored contains the stored model parameters.

Now we'll use this ARIMA model to predict new data points in the validation dataset using the PsiForecast() function. When array-entered into seven different Excel cells, this function will forecast six different future points in the dataset. (Note: The first forecasted point will be more accurate than the second, the second forecasted point more accurate than the third and so on.) The PsiForecast() function will be interactive in the sense that if any of the input values (values passed in the 2<sup>nd</sup> argument) change, the forecast will be recomputed.

The PsiForecastARIMA function takes five arguments but two are optional: Model, Input Data, Simulate, Num\_forecasts, and Header. Select a blank cell on the Data worksheet and enter =PsiForecast(. If using a version of Excel that does not support Dynamic Arrays, highlight cells B146:B152, then enter =PsiForecast(.

The first argument, Model, is the range of cells used by Analytic Solver Data Science to store the ARIMA model on the ARIMA\_Stored worksheet. This data range will change as the forecast method changes. Select or enter ARIMA\_Stored!B12:B38, for this argument.

The second argument, Input\_Data, is the range containing the initial starting points from the validation data set. The minimum number of initial points that should be specified for a seasonal ARIMA model is the larger of p + d + s \* (P +

D) and  $q + s * Q$ . In this example,  $p + d + s * (P + D)$  is equal to 13 ( $0 + 1 + 12 * (0 + 1)$ ) and  $q + s * Q$  is equal to 13 ( $1 + 12 * 1$ ), therefore the minimum number of initial starting points required is 13 (MAX (13, 13)). If you provide fewer than the minimum required number of starting points, PsiForecast() will return a column of zeros. (See the table below for the minimum number of initial starting points required by each Forecasting method included in Analytic Solver Data Science.) The maximum number of starting points is the number of points in the validation dataset. All points supplied in the second argument will be used in the forecast. Select or enter **Data!B1:B145**, for this argument.

Pass True or False for the third argument. Passing False will result in a static forecast that will only update if a cell passed in the 2<sup>nd</sup> argument is changed. If True is passed for this argument, a random error will be included in the forecasted points. See the Time Series Simulation example below for more information on passing True for this argument. In this case, Pass **False** for this argument.

Pass "7" for the next argument, number of forecasts.

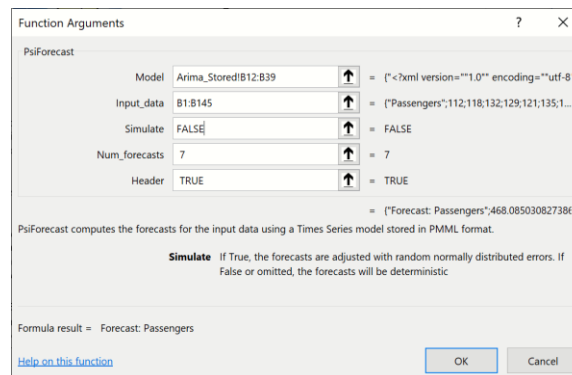
Pass "True" for header since to display a header at the top of the results.

Your formula should now be the following:

**=PsiForecast(ARIMA\_Stored!B12:B39,Data!B1:B145, False, 7, True).**

Note: If using a version of Excel that does not support Dynamic Arrays, press **CTRL + SHIFT + ENTER** to enter this formula as an array in all seven cells (B146:B152).

It's also possible to enter this formula using the Insert Function dialog by clicking Formulas – Insert Function, select PSI Data Science for Category, then PsiForecastARIMA.



The results from this function are displayed below.

Enter True for the Header argument to insert a heading above the forecasted values.

Notice that the formula is entered into cell B146 and the contents of the PsiForecast() Dynamic Array "spill" down into cells B147:B152.

|     | A      | B               | C |
|-----|--------|-----------------|---|
| 141 | Aug-60 | 606             |   |
| 142 | Sep-60 | 508             |   |
| 143 | Oct-60 | 461             |   |
| 144 | Nov-60 | 390             |   |
| 145 | Dec-60 | 432             |   |
| 146 |        | Forecast: 17899 |   |
| 147 |        | 468.082984      |   |
| 148 |        | 443.90313       |   |
| 149 |        | 444.106104      |   |
| 150 |        | 428.416324      |   |
| 151 |        | 424.563033      |   |
| 152 |        | 413.37003       |   |
| 153 |        |                 |   |

If any values change in the ranges ARIMA\_Stored!B12:B38 or Data!B2:B145, the forecast will be recomputed; but if the input argument values stay the same, the PsiForecast() function will always return the same forecast values. As mentioned above, the first forecasted value is the most accurate predicted point. Accuracy declines as the number of forecasted points increases.

See the section below for specifications on PsiForecast().

## PsiForecast()

---

`PsiForecast(Model, Input_Data, [Simulate], [Num_Forecasts], [Header])`

Computes the forecasts for Input\_Data using a Time Series model stored in PMML format.

**Model:** Range containing the stored Times Series model in PMML format.

**Input\_Data:** Range containing the new Time Series data for computing the forecasts. Range must contain a header with the time series name and a sufficient number of records for the forecasting with a given model.

**Simulate:** If True, the forecasts are adjusted with random normally distributed errors. If False or omitted, the forecasts will be deterministic.

**Num\_Forecasts:** Enter the number of desired forecasts.

**Header:** If True, a heading will be inserted above the forecasted results. If False or omitted, a heading will not be included in the result.

**Output:** A single column containing the header, if used, and forecasts for input time series. The number of produced forecasts is determined by the number of selected cells in the array-formula entry.

*Supported Models:*

- Arima
- Exponential Smoothing
- Double Exponential Smoothing
- Holt Winters Smoothing

*Previous related Psi Scoring functions:* PsiForecastARIMA, PsiForecastExp, PsiForecastDoubleExp, PsiForecastMovingAvg, PsiForecastHoltWinters

## Time Series Simulation

Analytic Solver Data Science includes the ability to perform a time series simulation, where future points in a time series are forecast on each Monte Carlo trial, using a model created via ARIMA or one of our smoothing methods (Exponential, Double Exponential, Moving Average, or Holt Winters).

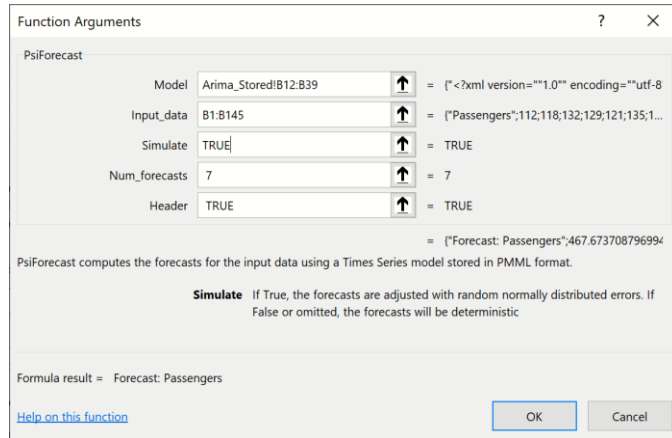
To run a time series simulation, we must pass “True” as the third argument to PsiForecast(). When the third argument is set to True, Analytic Solver will add a random (positive or negative) “epsilon” value to each forecasted point. Each time a simulation is run, 1000 trial “epsilon” values are generated using the PsiNormal distribution with parameters mean and standard deviation computed by the PsiForecast() function. You can view the output of this simulation in the same way as you would view “normal” simulation results in Analytic Solver Comprehensive, Analytic Solver Simulation or Analytic Solver Upgrade, simply by creating a PsiOutput() function and then double clicking the Output cell to view the Simulation Results dialog.

Select a blank cell, or Data!C146:C152 if using a version of Excel that does not support Dynamic Arrays, then click Formulas – Insert Function to display the Function Argument dialog.

As discussed previously, the first argument, ARIMA\_Stored!B12:B38, is the range of cells used by Analytic Solver to store the ARIMA model on the ARIMA\_Stored worksheet. This data range will change as the forecast method changes.

For the second argument the range containing the initial points in the series must be greater than the minimum number of initial points for a static forecast. For a seasonal ARIMA model when Simulate = True, the minimum number of initial points must be greater than  $\text{Max}((p + d + s * (P + D)), (q + s * Q))$ . In this example,  $p + d + s * (P + D)$  is equal to 13 ( $0 + 1 + 12 * (0 + 1)$ ) and  $q + s * Q$  is equal to 13 ( $1 + 12 * 1$ ), therefore the minimum number of initial starting points required is 13 (Minimum #Initial Points > MAX (13, 13)). However, when PsiForecastARIMA() is called with Simulate = True, it is recommended to add an additional number of datapoints, equal to the #Periods, to the minimum number required. In this instance the number of initial points will be 25: 13 (minimum # of points) + 12 (# of points for Period in the Time Series - ARIMA dialog). If you provide fewer than the minimum required number of starting points (13 in this example) PsiForecastARIMA() will return #VALUE. (See the table below for the minimum number of initial starting points required by each forecast method in Analytic Solver.) All points supplied in the second argument will be used in the forecast. Select or enter **Data!B1:B145**, for this argument.

Passing TRUE for the third argument indicates to Analytic Solver Data Science that you plan to use this function call in a Monte Carlo simulation, so it should add a random epsilon value (different on each Monte Carlo trial) to each forecasted point.



In versions of Excel supporting Dynamic Arrays, this formula is entered into cell C146 and the contents of the PsiForecast() Dynamic Array "spills" down into cells C147:C152. Note: If versions of Excel that do not support Dynamic Arrays, the formula must be entered as an Excel array.

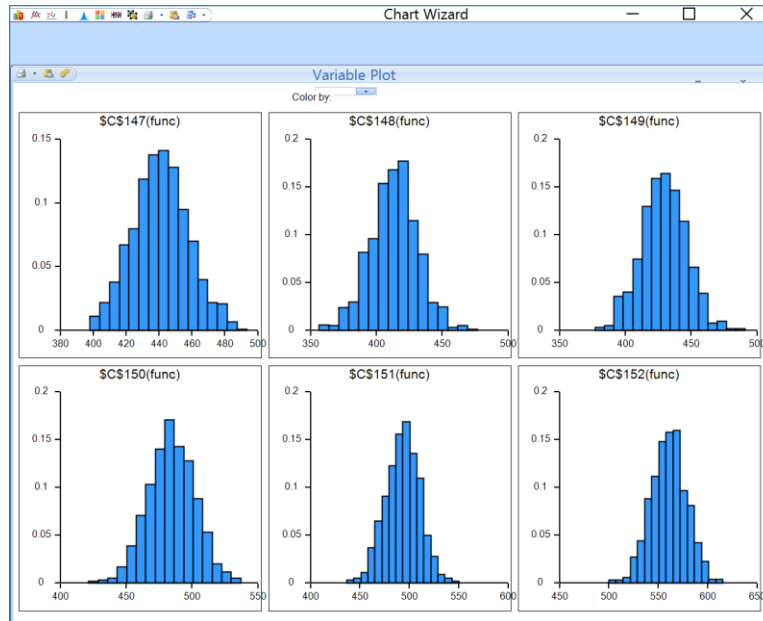
To view the results of the simulation including frequency and sensitivity charts, statistics, and percentiles for the full range of trial values in any version of Excel, we must first create an output cell. Select cell B147, then click Analytic Solver – Results – Referred Cell. Select cell D147 (or any blank cell on the spreadsheet) to enter the PsiOutput formula. Copy this formula from cell D147 down to cell D152. Therefore D147 = PsiOutput(C147), D148 = PsiOutput(C148), and so on.

|     | A      | B               | C               | D           |
|-----|--------|-----------------|-----------------|-------------|
| 143 | Oct-60 | 461             |                 |             |
| 144 | Nov-60 | 390             |                 |             |
| 145 | Dec-60 | 432             |                 |             |
| 146 |        | Forecast: Month | Forecast: Month |             |
| 147 | Jan-61 | 440.6349257     | 453.1571732     | 453.1571732 |
| 148 | Feb-61 | 412.1440457     | 409.1471791     | 409.1471791 |
| 149 | Mar-61 | 428.9072187     | 408.1621257     | 408.1621257 |
| 150 | Apr-61 | 485.2473609     | 496.8346369     | 496.8346369 |
| 151 | May-61 | 492.6189676     | 487.3958316     | 487.3958316 |
| 152 | Jun-61 | 560.7962281     | 576.0825778     | 576.0825778 |

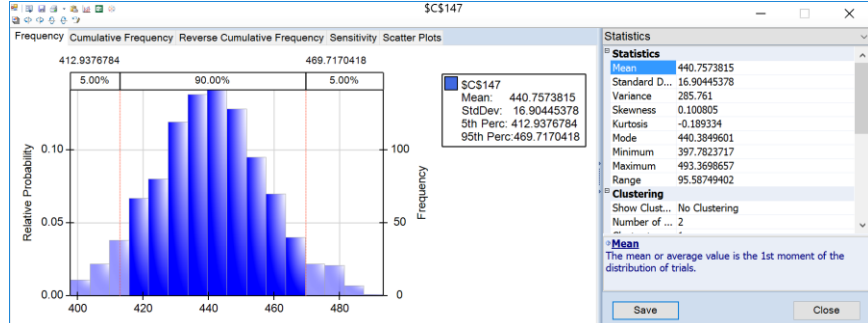
Click the down arrow on the Simulate icon and select Run Once. Instantly, Analytic Solver will perform a simulation with 1,000 Monte Carlo simulation trials (the default number). Since this is the first time a simulation has been performed, the following dialog opens. Subsequent simulations will not produce this report. However, it is possible to reopen the individual frequency charts by double clicking each of the output cells (B147:B152).

*Important Note: For Users who are familiar with simulation models in Analytic Solver Simulation, you'll notice that the time series simulation model that we just created now includes 6 uncertain functions, B147:B152, which are the cells containing our PsiForecast() functions. For more information on simulation with Analytic Solver, please see the Analytic Solver User Guide chapter, "Examples: Simulation and Risk Analysis".*

*PsiForecast() is not recognized as an uncertain function in the Cloud apps. If the simulation argument is set to "True", Analytic Solver App will generate a single random point around the forecast.*

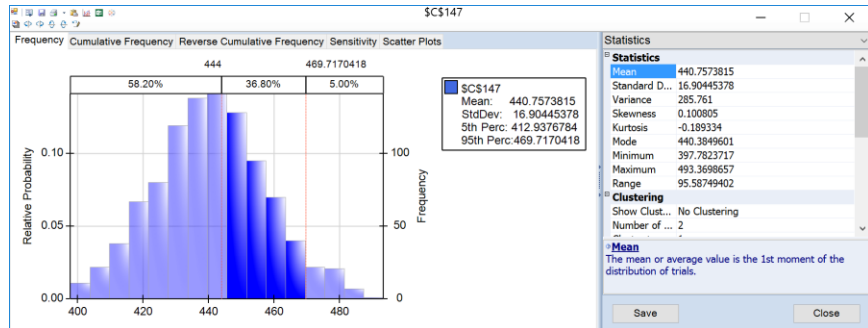


This dialog displays frequency charts for each of the six cells containing the forecasted data points. Double click the chart for cell C147 (top left) to open the Simulation Results dialog for the PsiForecast() function in cell C147. From here you can view frequency and sensitivity charts, statistics and percentiles for each forecasted point.



The frequency chart displays the distribution of all 1000 trial values for cell C147 with and observed mean 440.76 and standard deviation of 16.90 shown in the Chart Statistics. Select Simulate – Run Once a few more times (or click the green “play” button on the Solver Pane Model tab). Each time you do, another 1,000 Monte Carlo trials are run, and a slightly different mean will be displayed.

Enter 444 for the Lower Cutoff in the Chart Statistics section of the right panel, A vertical bar appears over the Frequency chart to display the frequency with which the forecasted value was greater than this value during the simulation. You can use this as an *estimate* of the probability that the actual value will be less than the forecasted value. In this case there was a 58.20% chance that the number of international airline passengers would be less than 444,000 in January 1961 and a 41.80% chance that the number of passengers would be greater than 444,000.



Looking to the right, you'll find the Statistics pane, which includes summary statistics for the full range of forecasted outcomes. We can see that the minimum forecasted value during this simulation was 397.78, and the maximum forecasted value was 493.37. Value at Risk 95% shows that 95% of the time, the number of international airline passengers was 469.72 or less in January 1961, in this simulation. The Conditional Value at Risk 95% value indicates that the average number of passengers we would have seen (up to the 95% percentile) was 438.86. For more information on Analytic Solver Platform's full range of features, see the Analytic Solver User Guide chapter, "Examples: Simulation and Risk Analysis".

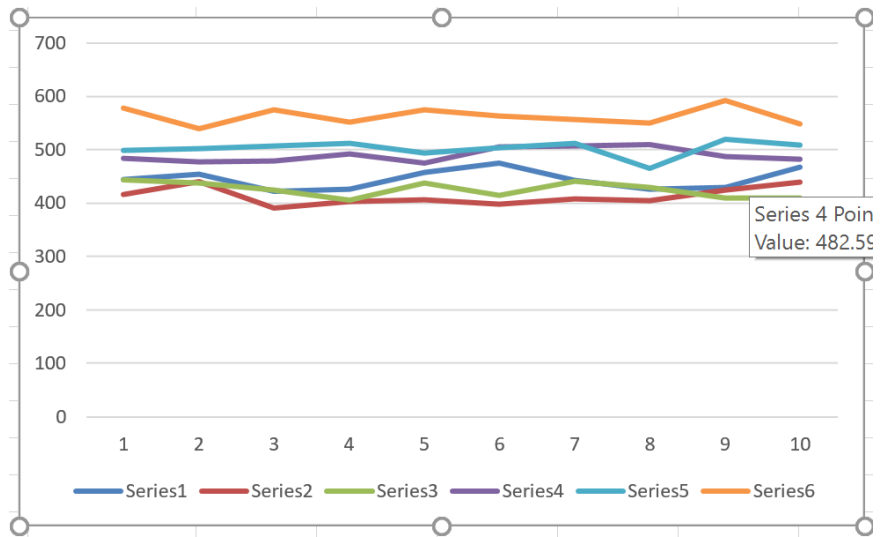
| Statistics                 |               |
|----------------------------|---------------|
| <b>Statistics</b>          |               |
| Mean                       | 440.7573815   |
| Standard Deviation         | 16.90445378   |
| Variance                   | 285.761       |
| Skewness                   | 0.100805      |
| Kurtosis                   | -0.189334     |
| Mode                       | 440.3849601   |
| Minimum                    | 397.7823717   |
| Maximum                    | 493.3698657   |
| Range                      | 95.58749402   |
| <b>Clustering</b>          |               |
| Show Clusters              | No Clustering |
| Number of Clusters         | 2             |
| Cluster to show            | 1             |
| <b>Advanced Statistics</b> |               |
| Mean Abs. Deviation        | 13.50426773   |
| SemiVariance               | 139.181       |
| SemiDeviation              | 11.7975       |
| Value at Risk 95%          | 469.7170418   |
| Cond. Value at Risk 95%    | 438.8635718   |
| Mean Confidence 95%        | 1.04773       |
| Std. Dev. Confidence 95%   | 0.70527       |
| Coefficient of Variation   | 0.0383532     |
| Standard Error             | 0.534298      |
| Expected Loss              | 0             |
| Expected Loss Ratio        | 0%            |
| Expected Gain              | 440.7573815   |
| Expected Gain Ratio        | 100%          |
| Expected Value Margin      | 1             |

Select cells E147:E156 and then enter the formula, =PsiData(C147), then press CTRL + SHIFT + ENTER to array enter the formula into all 10 cells. Repeat the same steps to array enter "=PsiData(C148)" in cells F147:F156, "=PsiData(C149)" in cells G147:G156, "=PsiData(C150)" in cells H147:H156,

“PsiData(C151)” in cells I147:I156, and “=PsiData(C152)” in cells J147:J156. Then click Simulate – Run Once to run a simulation.

The ten Excel cells in these columns will update with trial values for each of the PsiForecast() functions in column C. For example, cells E147:N147 will contain the first 10 trial values for the PsiForecast() function in cell C147, Cells E147:N147 will contain the first 10 trial values for cell C148 and so on. (For more information on the PsiData() function, please see the Excel Solvers Reference Guide chapter, “Psi Function Reference.”)

If we create an Excel chart of these values, you’ll see a chart similar to the one below where each of Series1 through Series6 represents a different Monte Carlo trial. The random “epsilon” value added to each forecast value accounts for (all of) the variation among the lines. If the third argument were FALSE or omitted, all of the lines would overlap, assuming that the table or parameters and the starting values were not changing.



The remaining Forecasting methods can be used in the same way using PsiForecast() with information from their respective Stored Model sheets.

| Forecasting Algorithm        | Stored Model Sheet    | Minimum # of Initial Points when Simulate = False | Minimum # of Initial Points when Simulate = True |
|------------------------------|-----------------------|---------------------------------------------------|--------------------------------------------------|
| Non- Seasonal ARIMA          | ARIMA_Stored          | Max(p + d, q)                                     | Max(p + d, q)                                    |
| Seasonal ARIMA               | ARIMA_Stored          | Max((p + d + s *(P + D), (q + s * Q)              | 1 + Max((p + d + s *(P + D), (q + s * Q)**       |
| Exponential Smoothing        | Expo_Stored           | 1                                                 | 1                                                |
| Double Exponential Smoothing | DoubleExpo_Stored     | 1                                                 | 1                                                |
| Moving Average Smoothing     | MovingAvg_Stored      | # of Intervals                                    | # of Intervals                                   |
| Holt Winters Smoothing       | MulHoltWinters_Stored | 2 * #Periods                                      | 2 * #Periods                                     |



|  |                           |  |  |
|--|---------------------------|--|--|
|  | AddHoltWinters_Stored     |  |  |
|  | NoTrendHoltWinters_Stored |  |  |

\*\*Adding a number of data points equal to the Number of Periods (as shown on the Time Series – ARIMA dialog) to the Minimum # of Initial Points when Simulate = True is recommended when calling PsiForecast() with Simulate = True.

---

## Scoring to a Database

This example describes the steps required to create a classification model using the Discriminant Analysis classification algorithm and then uses that model to score new data. *Note that this is only supported in Analytic Solver Desktop. This is not supported in the Data Science Cloud app.*

The example dataset Boston\_Housing.xlsx will be used to illustrate the steps required. Recall that this example dataset includes 14 variables related to housing prices collected from census tracts in the Boston area. For more information on this example dataset and Discriminant Analysis in general, please see the Discriminant Analysis chapter within the Analytic Solver Data Science Reference Guide.

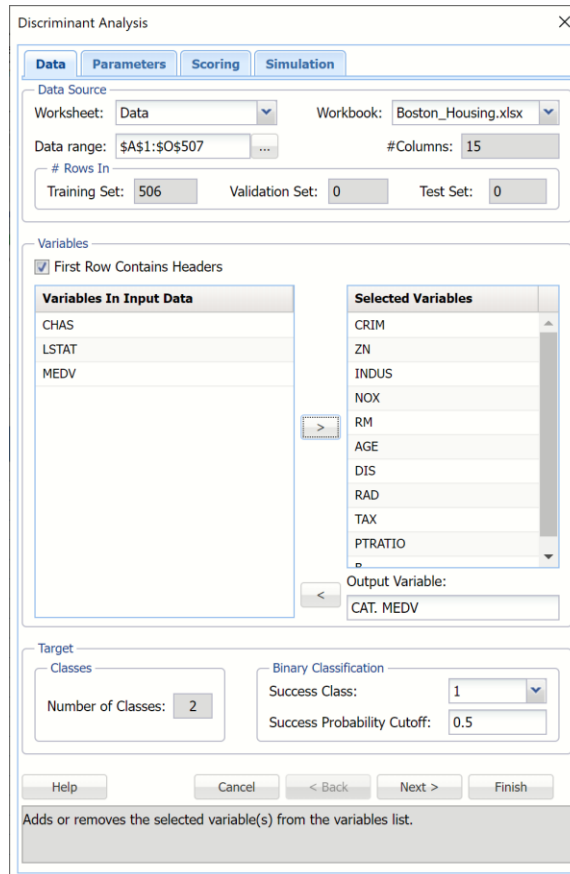
Open Boston\_Housing.xlsx, then click **Classify – Discriminant Analysis** to open the *Discriminant Analysis – Data* tab.

Select the **CAT.MEDV** variable in the *Variables In Input Data* list box then click > to select as the *Output Variable*. Immediately, the options for *Classes in the Output Variable* are enabled. *#Classes* is prefilled as “2” since the *CAT.MEDV* variable contains two classes, 0 and 1.

“*Success*” *Class (for Lift Chart)* is selected by default and Class 1 is to be considered a “success” or the significant class in the Lift Chart. (Note: This option is enabled when the number of classes in the output variable is equal to 2.)

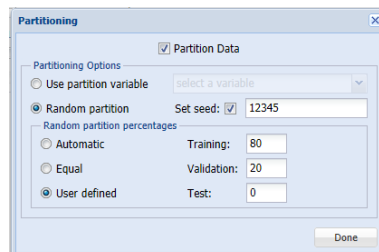
Enter a value between 0 and 1 here to denote the *Specify initial cutoff probability for success*. If the calculated probability for success for an observation is greater than or equal to this value, then a “success” (or a 1) will be predicted for that observation. If the calculated probability for success for an observation is less than this value, then a “non-success” (or a 0) will be predicted for that observation. The default value is 0.5. (Note: This option is only enabled when the # of classes is equal to 2.)

Select **CRIM, ZN, INDUS, NOX, RM, AGE, DIS, RAD, TAX, PTRATIO, & B** in the *Variables In Input Data* list box then click > to move to the *Selected Variables* list box. (*CHAS, LSTAT, & MEDV* should remain in the *Variables In Input Data* list box as shown below.)

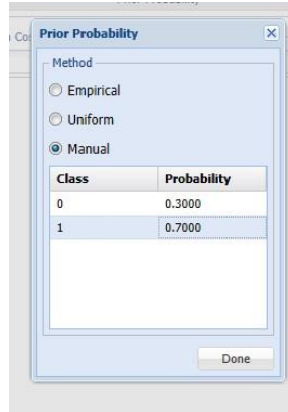


Click **Next** to advance to the Parameters tab.

Since we did not partition the dataset before we started the classification method, we can partition the dataset now. Click **Partition Data** and then select the Partition Data option to enable the Partitioning Options. Select **User Defined**, then enter **80** for *Training* and ensure **20** is automatically entered for *Validation*. Click Done to close the dialog.



Click Prior Probability to open the Prior Probability dialog. Three options appear: Empirical, Uniform, and Manual.

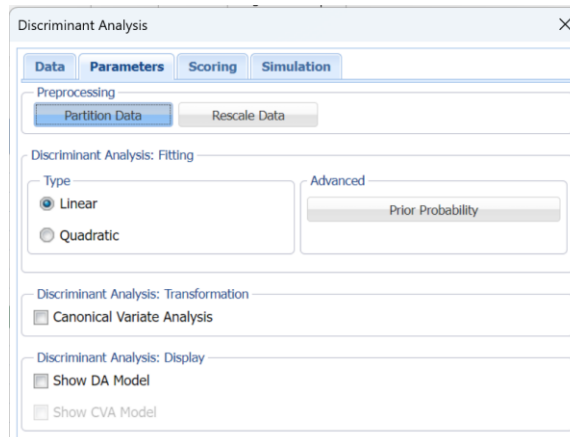


If the first option is selected, *Empirical*, Analytic Solver Data Science will assume that the probability of encountering a particular class in the dataset is the same as the frequency with which it occurs in the training data.

If the second option is selected, *Uniform*, Analytic Solver Data Science will assume that all classes occur with equal probability.

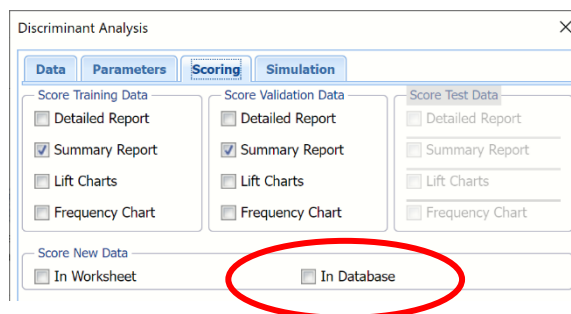
Select the third option, *Manual*, to manually enter the desired probability values for each class.

Under *Probability*, enter **0.7** for Class 1 and **0.3** for Class 0. Click Done to close the dialog.



Click **Next** to advance to the *Discriminant Analysis – Scoring* tab.

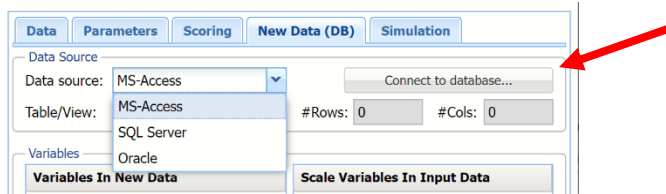
Since we did not create a test partition, the options for Score test data are disabled. See the chapter “Data Science Partitioning” within the Analytic Solver Data Science Reference Guide for information on how to create a test partition.



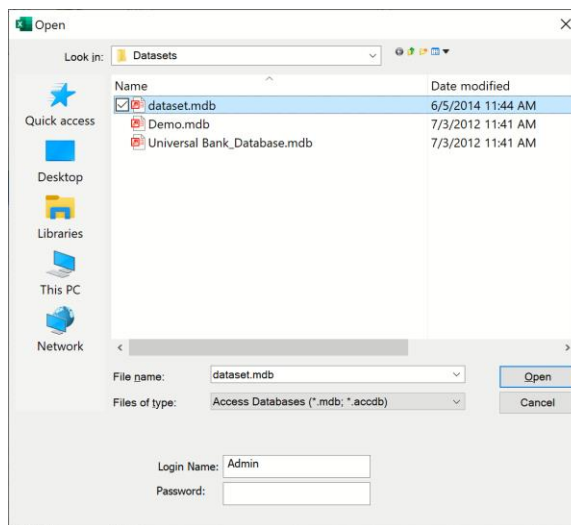
In the *Score New Data* group, select **In Database**. The *Scoring to Database* tab opens.

The first step on this tab is to select the **Data source**. Once the *Data source* is selected, **Connect to database...** will be enabled.

This example illustrates how to score to an MS-Access database. Select **MS-Access** for the *Data source*, then click **Connect to Database...**

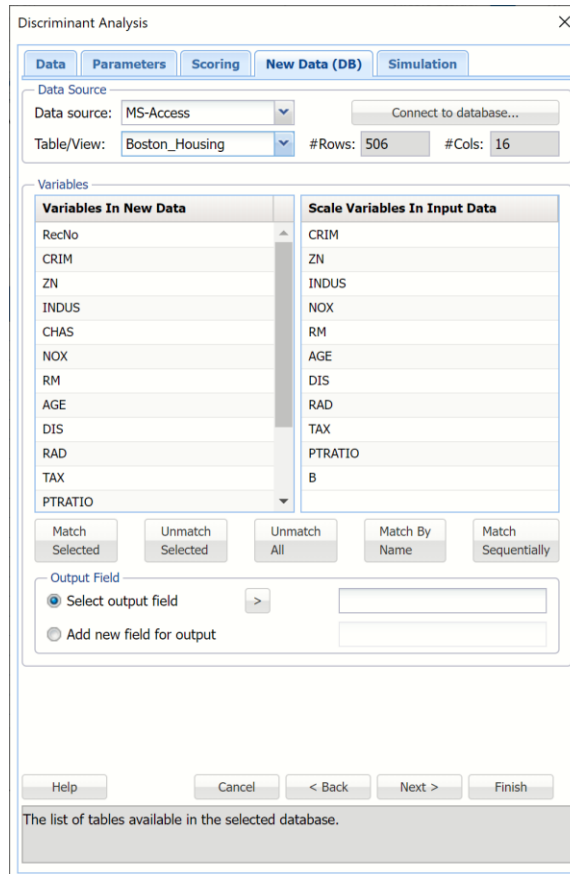


An *Open* dialog appears, browse to the location where dataset.mdb is saved, ...



...then click **Open**. Note the *Login Name* and *Password* fields at the bottom of the tab. If your database is password protected, enter the appropriate information here.

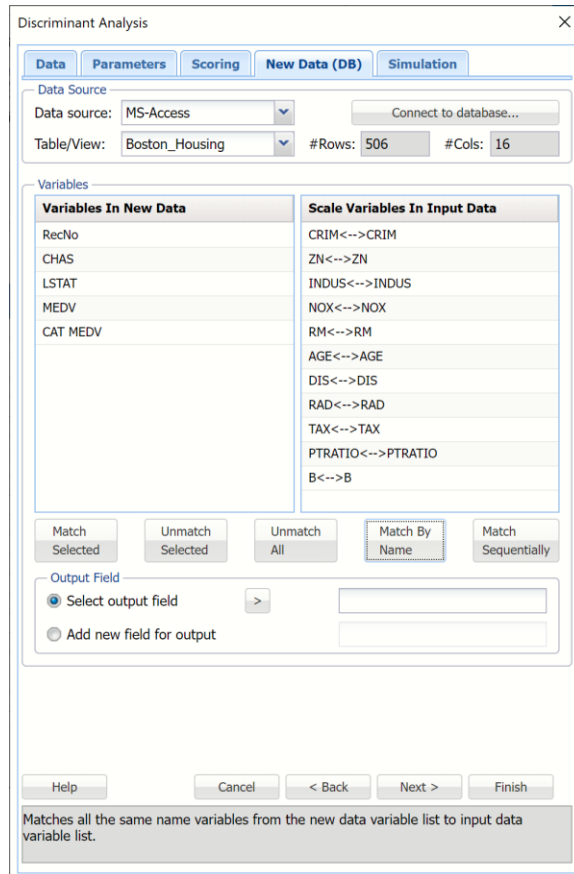
The *Scoring to Database* tab re-appears. Select **Boston\_Housing** for *Table/View*. The tab will be populated with variables from the database, dataset.mdb, under *Variables in New Data* and with variables from the Boston\_Housing.xlsx dataset under *Scale Variables In Input Data*.



Analytic Solver Data Science offers three easy techniques to match variables in the dataset to variables in the database:

1. Matching by Name.
2. Matching Sequentially
3. Manually Matching.

If **Match By Name** is clicked, all similar named variables in Boston\_Housing.xlsx will be matched with similar named variables in dataset.mdb, as shown in the screenshot below. Note that the additional database fields remain in the *Variables in New Data* list box while all variables in the *Scale Variables In Input Data* list box have been matched.



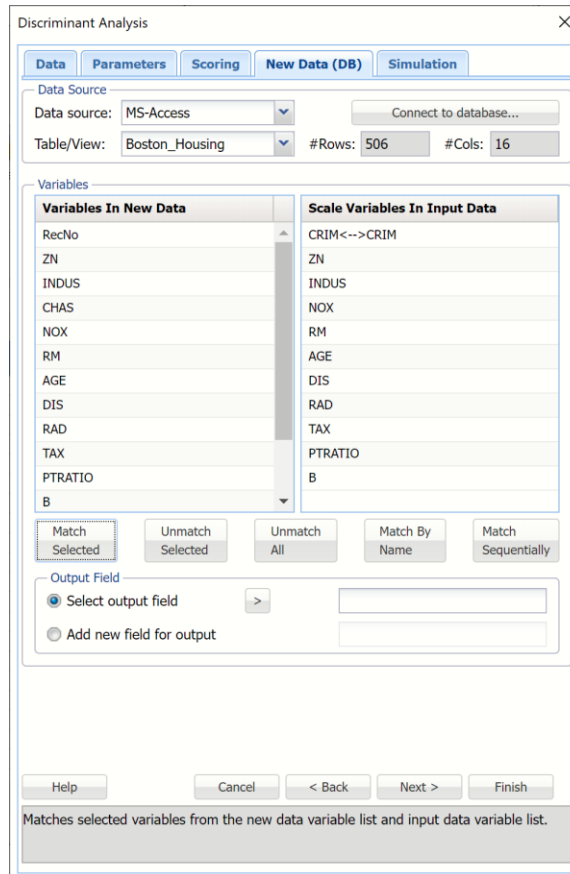
If **Match Sequentially** is clicked, the first 11 variables in Boston\_Housing.xlsx will be matched with the first 11 variables in the dataset.mdb database.

The first 11 variables in both the database and the dataset are now matched under *Scale Variables In Input Data*. The additional database fields remain under *Variables in New Data*.

Note: If Match Sequentially is clicked, variables from each list are matched in sequence, i.e. 1st variable under Variables in New Data is matched with the 1st variable under Scale Variables in Input Data, 2nd variable under Variables in New Data is matched with the 2nd variable under Scale Variables in Input Data, etc.

To manually map variables from the dataset to the database, select a field from the database in the *Variables in New Data* list box, then select the variable to be matched in the dataset in the *Scale Variables In Input Data* list box, then click **Match Selected**.

For example to match the CRIM variable in the database to the CRIM variable in the dataset, select **CRIM** from the dataset.mdb database in the *Variables in New Data* list box, select **CRIM** from *Scale Variables In Input Data*, then click **Match Selected** to match the two variables.



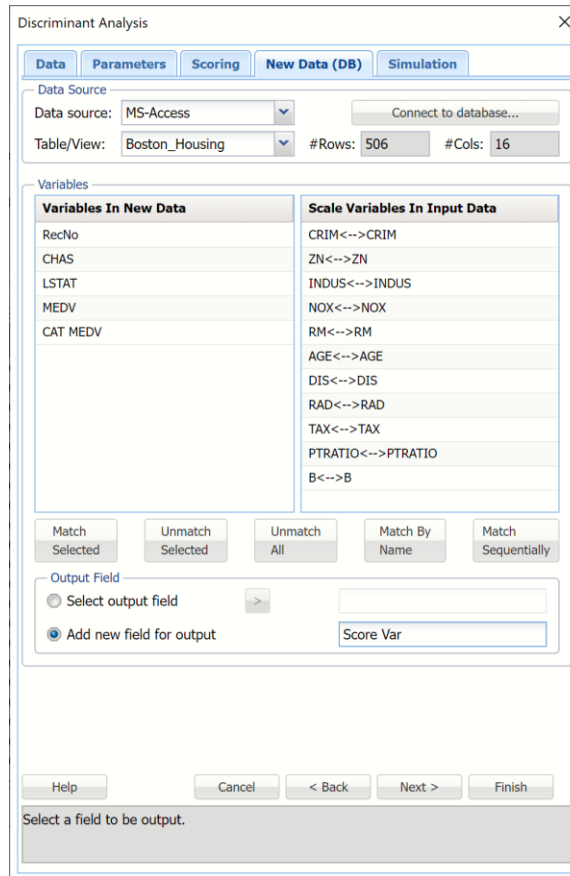
Notice that *CRIM* has been removed from the *Variables in New Data* list box and is now listed next to *CRIM* in the *Scale Variables In input data* list box. Continue with these steps to match the remaining 10 variables in the *Boston\_Housing.xlsx* dataset.

To unmatch all variables click **Unmatch all**. To unmatch a single match, highlight the match, then click **Unmatch selected**.

An Output Field may be selected from the remaining database fields listed under *Variables in New Data* or a new Output Field can be added. *Note: An output field must be a string.*

For this example,

1. Click **Unmatch All**,
2. Click **Match By Name**,
3. Select **Add new field for output**,
4. Enter a name in the field to the right such as “Score Var”.



From here you can either click **Finish** to score the new data or you can click **Next** to advance to the **Simulation** tab where you can perform a complete risk analysis on the data. See the previous chapter, *Automatic Risk Analysis of Machine Learning Models*, help topic for more information on the options displayed on the **Simulation** tab.

Click **Finish** on the *New Data (DB)* tab. The *DA\_NewScoreDB* worksheet is inserted to the right of the *Data* worksheet. In addition, the worksheet name is also inserted into the **Model** tab of the **Analytic Solver** task pane under **Reports – Discriminant Analysis**. This worksheet simply includes the name of the database, Table Name, # of records scored and the variables.

|    | B                    | C                                               | D  | E     | F   | G  | H   | I   | J   | K   | L       | M | N |
|----|----------------------|-------------------------------------------------|----|-------|-----|----|-----|-----|-----|-----|---------|---|---|
| 10 | <b>Inputs</b>        |                                                 |    |       |     |    |     |     |     |     |         |   |   |
| 11 |                      |                                                 |    |       |     |    |     |     |     |     |         |   |   |
| 12 | <b>Data</b>          |                                                 |    |       |     |    |     |     |     |     |         |   |   |
| 13 | Database Name        | C:\Users\Nicole\Documents\Frontline\dataset.mdb |    |       |     |    |     |     |     |     |         |   |   |
| 14 | Table Name           | Boston_Housing                                  |    |       |     |    |     |     |     |     |         |   |   |
| 15 | # Records Scored     | 506                                             |    |       |     |    |     |     |     |     |         |   |   |
| 16 |                      |                                                 |    |       |     |    |     |     |     |     |         |   |   |
| 17 | <b>Variables</b>     |                                                 |    |       |     |    |     |     |     |     |         |   |   |
| 18 | # Input Variables    | 11                                              |    |       |     |    |     |     |     |     |         |   |   |
| 19 | Input variables      | CRIM                                            | ZN | INDUS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B |   |
| 20 | Corresponding Fields | CRIM                                            | ZN | INDUS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B |   |
| 21 | Output Field         | Score Var                                       |    |       |     |    |     |     |     |     |         |   |   |
| 22 |                      |                                                 |    |       |     |    |     |     |     |     |         |   |   |

To view the scored records, open the *dataset.mdb* database in Microsoft Access and inspect the *Score Var* column as shown in the screenshot below. (Click the **Enable Content** button to view the dataset.)

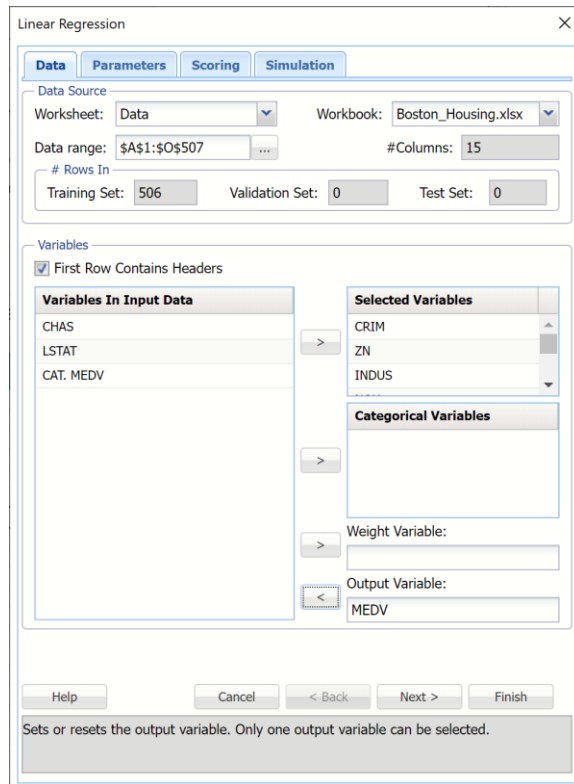


The screenshot shows a software interface with a 'Tables' pane on the left and a data table for 'Boston\_Housing' on the right. The 'Tables' pane lists various tables including 'Assoc\_binary', 'Assoc\_Itemist', 'Boston\_Housing', 'Charles\_BookClub', 'Digits', 'Flying\_Fitness', 'Iris', 'Utilities', and 'Wine'. The 'Boston\_Housing' table is selected and its data is displayed in a grid. The grid has columns for 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV', 'CAT MEDV', and 'Score Var'. The first row of data shows values: 0, 0.538, 6.575, 65.2, 4.09, 1, 296, 15.3, 396.9, 4.98, 24, 0, 1. A red arrow points to the top right corner of the data table.

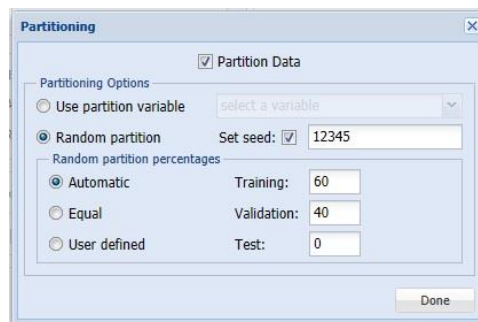
Note: If the database has been saved to a location that does not allow write-access, you will get the error, [\[Microsoft\]\[ODBC Microsoft Access Driver\] Cannot modify the design of table 'Boston\\_Housing'. It is in a read-only database.](#)

## Scoring to a Worksheet

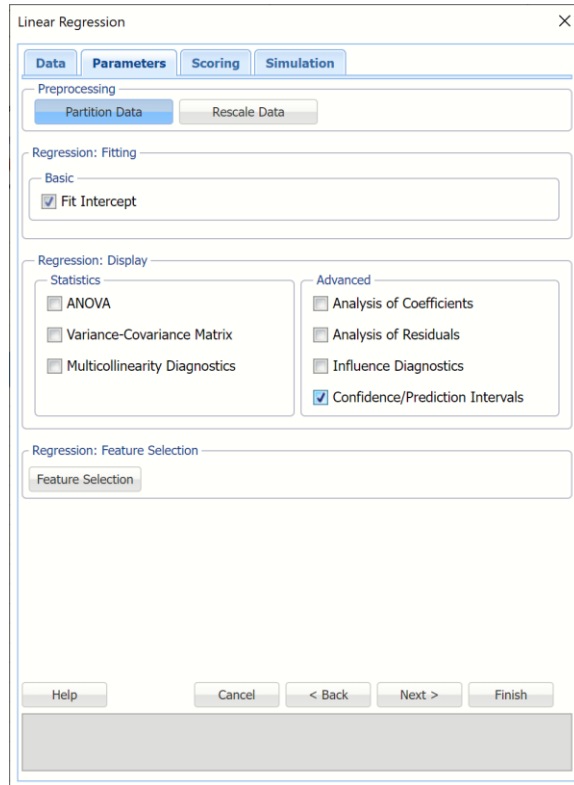
Analytic Solver Data Science can also perform scoring on new data in a worksheet. To illustrate, we'll re-use the Boston Housing example dataset. (To open, click Help – Example Models – Forecasting / Data Science Examples.) Click **Predict – Linear Regression** to open the *Linear Regression - Data* tab and select the Output and Selected Variables as shown in the screenshot below.



Click **Next** to advance to the Parameters tab. Since we haven't yet partitioned the dataset, we will do so now. Click Partition Data to open the Partition Data Dialog, then select the Partition Data option to enable the Standard Partitioning Options. For more information on these partitioning options, please see the Data Science Partitioning chapter that occurs earlier in this guide. Click Done to accept the partitioning defaults.

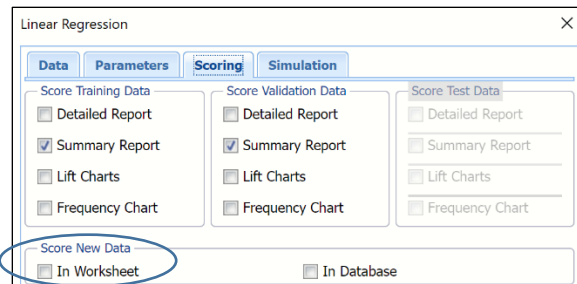


Under Advanced, Select Confidence/Prediction Intervals to include these intervals for the new predictions.

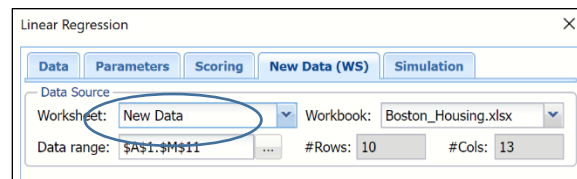


Click Next to advance to the Scoring tab.

Select **In Worksheet** in the *Score new data* group. The tab for *Match variables in the New Range* appears.



Select **New Data** for *Worksheet* at the top of the tab.



The variables listed under *Variables in New Data* are from the New Data worksheet and the variables listed under *Scale Variables In Input Data* are from the Data worksheet. Variables can be matched in three different ways.

1. Matching by Name.
2. Matching Sequentially
3. Manually Matching.

If *Match By Name* is clicked, variables with the same names in each set will be matched.

If *Match Sequentially* is clicked, variables from each list are matched in sequence, i.e. 1<sup>st</sup> variable under Variables in New Data is matched with the 1<sup>st</sup> variable under Scale Variables in Input Data, 2<sup>nd</sup> variable under Variables in New Data is matched with the 2<sup>nd</sup> variable under Scale Variables in Input Data, etc.

Variables may also be matched manually by selecting a variable under *Variables in New Data*, selecting a variable in *Scale Variables In Input Data*, and clicking *Match Selected*. For example, select **CRIM** under *Variables in new data* and **CRIM** under *Scale Variables In input data*, then click **Match Selected**.

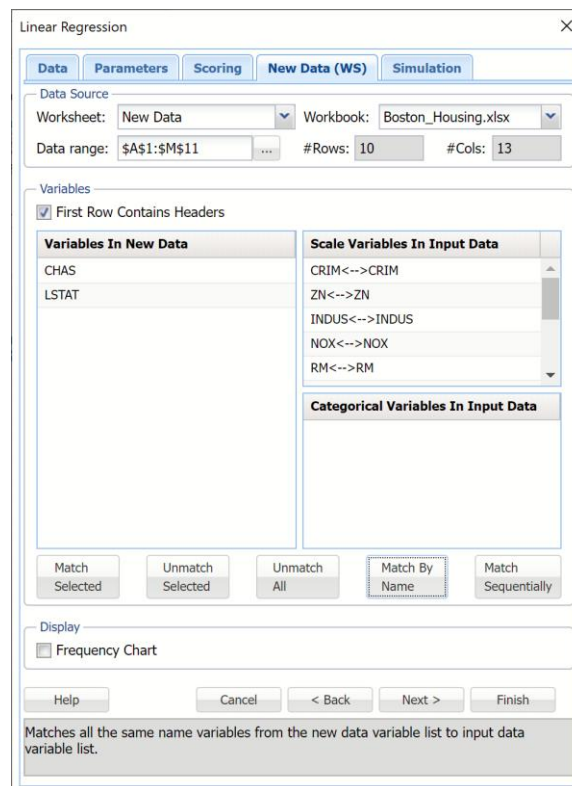
Click Match By Name to match the variables.

To unmatch all matched variables, click **Unmatch all**. To unmatch only one set of matched variables, select the matched variables in the *Scale Variables In input data* list box, then select **Unmatch Selected**.

Risk Analysis Options To unmatch only one set of matched variables, select the matched variables in the *Scale Variables In input data* list box, then select **Unmatch Selected**.

To add frequency charts to the output, select Frequency Chart. See the first section in the chapter, above, for the output from this checkbox. Leave this checkbox unchecked for this example.

From here you can either click Finish to score the new data or you can click Next to advance to the Simulation tab where you can perform a complete risk analysis on the data. See the previous chapter, Automatic Risk Analysis of Machine Learning Models, help topic for more information on the options displayed on the Simulation tab.



Click **Finish**.

Click the *LinReg\_NewScore* worksheet, inserted to the right of the Data tab, to view the output as shown below.

|    | B                              | C         | D                |
|----|--------------------------------|-----------|------------------|
| 10 | <b>New: Prediction Details</b> |           |                  |
| 11 |                                |           |                  |
| 12 |                                | Record ID | Prediction: MEDV |
| 13 |                                | Record 1  | 18.91758817      |
| 14 |                                | Record 2  | 15.69843282      |
| 15 |                                | Record 3  | 37.61687219      |
| 16 |                                | Record 4  | 8.573113358      |
| 17 |                                | Record 5  | 21.85853659      |
| 18 |                                | Record 6  | 14.44014655      |
| 19 |                                | Record 7  | 5.837839083      |
| 20 |                                | Record 8  | 49.69927467      |
| 21 |                                | Record 9  | -1.293189655     |
| 22 |                                | Record 10 | 3.481161659      |
| 23 |                                |           |                  |

Click the **Intervals: New** link on the Output Navigator to view the Confidence and Prediction Intervals on the *LinReg\_Intervals* worksheet, for the new predictions.

|    | P                     | Q         | R                     | S                     | T                     | U                     |
|----|-----------------------|-----------|-----------------------|-----------------------|-----------------------|-----------------------|
| 10 | <b>Intervals: New</b> |           |                       |                       |                       |                       |
| 11 |                       |           |                       |                       |                       |                       |
| 12 |                       | Record ID | 95% Confidence: Lower | 95% Confidence: Upper | 95% Prediction: Lower | 95% Prediction: Upper |
| 13 |                       | Record 1  | 16.2879985            | 21.54717784           | 7.174358385           | 30.86081794           |
| 14 |                       | Record 2  | 14.02418984           | 17.37287601           | 4.131592551           | 27.2852731            |
| 15 |                       | Record 3  | 31.84394081           | 43.58980357           | 24.70700724           | 50.52673713           |
| 16 |                       | Record 4  | 4.012379662           | 13.13384705           | -3.747154314          | 20.89338103           |
| 17 |                       | Record 5  | 17.56849368           | 26.14857949           | 9.63588664            | 34.08118654           |
| 18 |                       | Record 6  | 10.50283077           | 18.37746233           | 2.336793328           | 26.54349977           |
| 19 |                       | Record 7  | 1.758005139           | 9.916673028           | -6.312288076          | 17.98796624           |
| 20 |                       | Record 8  | 44.92337633           | 54.47517301           | 37.29774619           | 62.10080315           |
| 21 |                       | Record 9  | -10.67688564          | 8.090506328           | -16.09327256          | 13.50689325           |
| 22 |                       | Record 10 | -1.804371518          | 8.766694836           | -9.125406654          | 16.08772997           |
| 23 |                       |           |                       |                       |                       |                       |

Here we see the 95% Confidence and Prediction Intervals. Typically, Prediction Intervals are more widely utilized as they are a more robust range for the predicted value. For a given record, the Confidence Interval gives the mean value estimation with 95% probability. This means that with 95% probability, the regression line will pass through this interval. The Prediction Interval takes into account possible future deviations of the predicted response from the mean. There is a 95% chance that the predicted value will lie within the Prediction interval.

For more information on the rest of the Linear Regression output, please see the Predicting Housing Prices using Multiple Linear Regression chapter that appears earlier in this guide.